

---

# Appendix to “Auxiliary Task Reweighting for Minimum-data Learning”

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Additional Discussion on ARML

2 In this section we add more discussion on validity and soundness of ARML, especially on the three  
3 problems (**True Prior (P1)**, **Samples (P2)**, **Partition Function (P3)**), and how we resolve them  
4 (Sec. 3.3).

### 5 1.1 Full Version and Proof of Theorem 1 (P1)

6 In **True Prior (P1)** (Sec. 3.3) we use

$$\min_{\alpha} E_{\theta \sim p^J} \log \frac{p^m(\theta)}{p_{\alpha}(\theta)} \quad (\text{A1})$$

7 as a surrogate objective for the original optimization problem

$$\min_{\alpha} D_{\text{KL}}(p^*(\theta) \parallel p_{\alpha}(\theta)). \quad (\text{A2})$$

8 In this section, we will first intuitively explain why optimizing (A1) can end up with a near-optimal  
9 solution for (A2), and what assumptions do we need to make. Then we will give the full version of  
10 Theorem 1 and also the proof.

11 Let  $f(\alpha) = E_{\theta \sim p^J} \log \frac{p^m(\theta)}{p_{\alpha}(\theta)} = \frac{1}{Z(\alpha)} \int p^m(\theta) p_{\alpha}(\theta) \log \frac{p^m(\theta)}{p_{\alpha}(\theta)} d\theta$  be the optimization objective  
12 in (A1), where  $p^J(\theta) = \frac{p^m(\theta) p_{\alpha}(\theta)}{Z(\alpha)}$  and  $Z(\alpha) = \int p^m(\theta) p_{\alpha}(\theta) d\theta$  is the normalization term. Assume  
13  $p^*(\theta)$  has a compact support set  $S$ . Then we can write  $f(\alpha)$  as

$$\begin{aligned} f(\alpha) &= \frac{1}{Z(\alpha)} \int_{\theta \in S} p^m(\theta) p_{\alpha}(\theta) \log \frac{p^m(\theta)}{p_{\alpha}(\theta)} d\theta + \frac{1}{Z(\alpha)} \int_{\theta \notin S} p^m(\theta) p_{\alpha}(\theta) \log \frac{p^m(\theta)}{p_{\alpha}(\theta)} d\theta \\ &= \frac{Z(S; \alpha)}{Z(S; \alpha) + Z(\bar{S}; \alpha)} \int_{\theta \in S} \frac{p^m(\theta) p_{\alpha}(\theta)}{Z(S; \alpha)} \log \frac{p^m(\theta)}{p_{\alpha}(\theta)} d\theta \\ &\quad + \frac{Z(\bar{S}; \alpha)}{Z(S; \alpha) + Z(\bar{S}; \alpha)} \int_{\theta \notin S} \frac{p^m(\theta) p_{\alpha}(\theta)}{Z(\bar{S}; \alpha)} \log \frac{p^m(\theta)}{p_{\alpha}(\theta)} d\theta \\ &= f(\alpha; S) + f(\alpha; \bar{S}), \end{aligned} \quad (\text{A3})$$

14 where we denote the first and second term by  $f(\alpha; S)$  and  $f(\alpha; \bar{S})$  respectively,  $Z(S; \alpha) =$   
15  $\int_{\theta \in S} p^m(\theta) p_{\alpha}(\theta) d\theta$  and  $Z(\bar{S}; \alpha) = \int_{\theta \notin S} p^m(\theta) p_{\alpha}(\theta) d\theta$  are the normalization terms inside and  
16 outside  $S$ .

17 To build the connection between the surrogate objective  $f(\alpha)$  and the original objective  $KL_{\alpha} :=$   
18  $D_{\text{KL}}(p^*(\theta) \parallel p_{\alpha}(\theta))$ , we make the following assumption,

19 **Assumption 1.** *The support set  $S$  is small so that  $p_{\alpha}(\theta)$  and  $p^m(\theta)$  are constants inside  $S$ , and  $p^*(\theta)$   
20 is uniform in  $S$ .*

21 This assumption is reasonable when  $S$  is really informative, which we assume is the case for the true  
 22 prior  $p^*(\theta)$  [3]. With this assumption, we have

$$KL_\alpha = \int_{\theta \in S} p^*(\theta) \log \frac{p^*(\theta)}{p_\alpha(\theta)} d\theta = \log \frac{p^*(\theta^*)}{p_\alpha(\theta^*)} \cdot \int_{\theta \in S} p^*(\theta) d\theta = \log \frac{p^*(\theta^*)}{p_\alpha(\theta^*)}, \quad (\text{A4})$$

23 where  $\theta^* \in S$  is the optimal parameter. We can also write  $f(\alpha; S)$  as

$$\begin{aligned} f(\alpha; S) &= \frac{Z(S; \alpha)}{Z(S; \alpha) + Z(\bar{S}; \alpha)} \int_{\theta \in S} \frac{p^m(\theta) p_\alpha(\theta)}{Z(S; \alpha)} \log \frac{p^m(\theta)}{p_\alpha(\theta)} d\theta \\ &= \frac{Z(S; \alpha)}{Z(S; \alpha) + Z(\bar{S}; \alpha)} \log \frac{p^m(\theta^*)}{p_\alpha(\theta^*)} \cdot \int_{\theta \in S} \frac{p^m(\theta) p_\alpha(\theta)}{Z(S; \alpha)} d\theta \\ &= \frac{Z(S; \alpha)}{Z(S; \alpha) + Z(\bar{S}; \alpha)} \log \frac{p^m(\theta^*)}{p_\alpha(\theta^*)} \\ &= \frac{Z(S; \alpha)}{Z(S; \alpha) + Z(\bar{S}; \alpha)} (\log \frac{p^*(\theta^*)}{p_\alpha(\theta^*)} + \log \frac{p^m(\theta^*)}{p^*(\theta^*)}) \\ &= \frac{Z(S; \alpha)}{Z(S; \alpha) + Z(\bar{S}; \alpha)} (KL_\alpha + C_1), \end{aligned} \quad (\text{A5})$$

24 where  $C_1 = \log \frac{p^m(\theta^*)}{p^*(\theta^*)}$  is a constant invariant to  $\alpha$ . Since  $p^m(\theta)$  also covers other ‘‘overfitting’’ area  
 25 other than  $S$ , we can assume that  $p^*(\theta^*) \geq p^m(\theta^*)$ , which gives  $C_1 \leq 0$ . Furthermore, we can notice  
 26 that

$$Z(S; \alpha) = \int_{\theta \in S} p^m(\theta) p_\alpha(\theta) d\theta = \int_{\theta \in S} \frac{p^m(\theta) p_\alpha(\theta)}{p^*(\theta)} p^*(\theta) d\theta = \frac{p^m(\theta^*) p_\alpha(\theta^*)}{p^*(\theta^*)} = C_2 e^{-KL_\alpha}, \quad (\text{A6})$$

27 where  $C_2 = p^m(\theta^*)$  is a constant invariant to  $\alpha$ . Then we can write  $f(\alpha; S)$  as

$$f(\alpha; S) = \frac{C_2 e^{-KL_\alpha}}{C_2 e^{-KL_\alpha} + Z(\bar{S}; \alpha)} (KL_\alpha + C_1). \quad (\text{A7})$$

28 In this way, we build the connection between the surrogate objective  $f(\alpha)$  and the original objective  
 29  $KL_\alpha$ .

30 Now we give an intuitive explanation for why optimizing  $f(\alpha)$  gives a small  $KL_\alpha$  as well. We can  
 31 write  $f(\alpha)$  as

$$\begin{aligned} f(\alpha) &= f(\alpha; S) + f(\alpha; \bar{S}) \\ &= \frac{C_2 e^{-KL_\alpha}}{C_2 e^{-KL_\alpha} + Z(\bar{S}; \alpha)} (KL_\alpha + C_1) + \frac{Z(\bar{S}; \alpha)}{C_2 e^{-KL_\alpha} + Z(\bar{S}; \alpha)} \int_{\theta \in \bar{S}} \frac{p^m(\theta) p_\alpha(\theta)}{Z(\bar{S}; \alpha)} \log \frac{p^m(\theta)}{p_\alpha(\theta)} d\theta. \end{aligned} \quad (\text{A8})$$

32 As one can notice,  $f(\alpha)$  not only depends on  $KL_\alpha$ , but also on  $Z(\bar{S}; \alpha)$  and the integral  
 33  $\int_{\theta \in \bar{S}} \frac{p^m(\theta) p_\alpha(\theta)}{Z(\bar{S}; \alpha)} \log \frac{p^m(\theta)}{p_\alpha(\theta)} d\theta$ . First we remove the dependency on the integral by taking its lower  
 34 bound and upper bound. Concretely, with Jensen’s inequality, we have

$$\int_{\theta \in \bar{S}} \frac{p^m(\theta) p_\alpha(\theta)}{Z(\bar{S}; \alpha)} \log \frac{p^m(\theta)}{p_\alpha(\theta)} d\theta \leq \log \frac{\int_{\theta \in \bar{S}} (p^m(\theta))^2 d\theta}{Z(\bar{S}; \alpha)} = \log \frac{C_3}{Z(\bar{S}; \alpha)}, \quad (\text{A9})$$

35 where  $C_3 = \int_{\theta \in \bar{S}} (p^m(\theta))^2 d\theta$  is a constant invariant to  $\alpha$ . Likewise, we have

$$\begin{aligned} \int_{\theta \in \bar{S}} \frac{p^m(\theta) p_\alpha(\theta)}{Z(\bar{S}; \alpha)} \log \frac{p^m(\theta)}{p_\alpha(\theta)} d\theta &= \int_{\theta \in \bar{S}} -\frac{p^m(\theta) p_\alpha(\theta)}{Z(\bar{S}; \alpha)} \log \frac{p_\alpha(\theta)}{p^m(\theta)} d\theta \\ &\geq -\log \frac{\int_{\theta \in \bar{S}} (p_\alpha(\theta))^2 d\theta}{Z(\bar{S}; \alpha)} \\ &\geq -\log \frac{C_4}{Z(\bar{S}; \alpha)}, \end{aligned} \quad (\text{A10})$$

36 where  $C_4 = \max_\alpha \int_{\theta \in \bar{S}} (p_\alpha(\theta))^2 d\theta$  is a constant invariant to  $\alpha$ . In this way, we get the lower bound  
 37 and upper bound for  $f(\alpha)$ :

$$\begin{aligned} f(\alpha) &\geq f_l(\alpha) = \frac{C_2 e^{-KL_\alpha}}{C_2 e^{-KL_\alpha} + Z(\bar{S}; \alpha)} (KL_\alpha + C_1) - \frac{Z(\bar{S}; \alpha)}{C_2 e^{-KL_\alpha} + Z(\bar{S}; \alpha)} \log \frac{C_4}{Z(\bar{S}; \alpha)}, \\ f(\alpha) &\leq f_u(\alpha) = \frac{C_2 e^{-KL_\alpha}}{C_2 e^{-KL_\alpha} + Z(\bar{S}; \alpha)} (KL_\alpha + C_1) + \frac{Z(\bar{S}; \alpha)}{C_2 e^{-KL_\alpha} + Z(\bar{S}; \alpha)} \log \frac{C_3}{Z(\bar{S}; \alpha)}. \end{aligned} \quad (\text{A11})$$

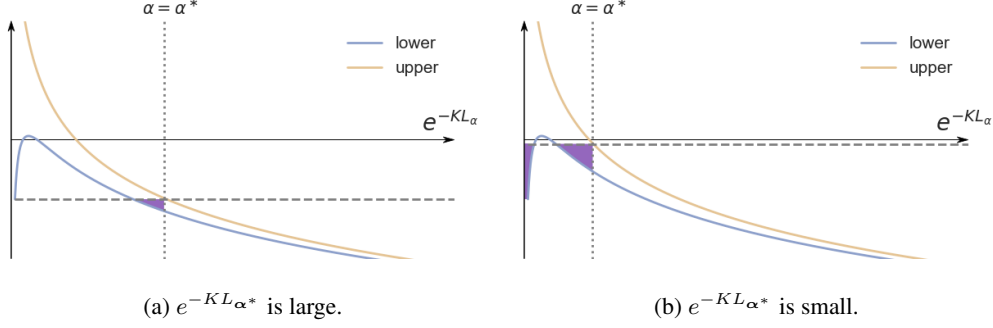


Figure 1:  $f(e^{-KL\alpha})$ 's upper bound  $f_u(e^{-KL\alpha})$  (golden line) and lower bound  $f_l(e^{-KL\alpha})$  (blue line).  $\alpha^* = \arg \max_{\alpha} (e^{-KL\alpha}) = \arg \min_{\alpha} KL_{\alpha}$  denotes the largest  $e^{-KL\alpha}$  we could possibly reach. Shaded region denotes where  $(e^{-KL\hat{\alpha}}, f(e^{-KL\hat{\alpha}}))$  could possibly be.

38 We plot  $f_l$  and  $f_u$  as functions of  $e^{-KL\alpha}$  in Fig. 1 (here we assume  $Z(\bar{S}; \alpha)$  is constant w.r.t.  $\alpha$  for  
 39 brevity).  $f(\alpha)$  lies between the upper bound (golden line) and the lower bound (blue line).

40 Our goal is to find the optimal  $\alpha^*$  that minimizes  $KL_{\alpha}$ , i.e.,  $\alpha^* = \arg \min_{\alpha} KL_{\alpha} =$   
 41  $\arg \max_{\alpha} e^{-KL\alpha}$ . By optimizing  $f(\alpha)$ , we end up with a suboptimal  $\hat{\alpha} = \arg \min_{\alpha} f(\alpha)$ . Ideally,  
 42 we hope that  $KL_{\hat{\alpha}}$  is close to  $KL_{\alpha^*}$ , which means when we minimize  $f(\hat{\alpha})$ , we can also get a large  
 43  $e^{-KL\hat{\alpha}}$ . This is the case when  $e^{-KL\alpha^*}$  is large (see Fig. 1a). When  $e^{-KL\alpha^*}$  is large, the upper  
 44 bound  $f_u$  and the lower bound  $f_l$  are close to each other around  $e^{-KL\alpha^*}$  (this is the case when  
 45  $Z(\bar{S}; \alpha)$  is small). Since we have

$$f_l(e^{-KL\hat{\alpha}}) \leq f(e^{-KL\hat{\alpha}}) \leq f(e^{-KL\alpha^*}) \leq f_u(e^{-KL\alpha^*}), \quad (\text{A12})$$

46 we can assert that  $(e^{-KL\hat{\alpha}}, f(e^{-KL\hat{\alpha}}))$  lies in the shaded region, because if  $e^{-KL\hat{\alpha}}$  is on the left side  
 47 of the region, we have  $f(e^{-KL\hat{\alpha}}) \geq f_u(e^{-KL\alpha^*})$  which is contradictory to (A12), and if  $e^{-KL\hat{\alpha}}$   
 48 cannot be on the right side of the region because  $e^{-KL\alpha^*}$  is the furthest we can go. Since the shaded  
 49 region is small,  $KL_{\hat{\alpha}}$  is thus close to the optimal solution  $KL_{\alpha^*}$ .

50 Unfortunately, this may not hold anymore when  $e^{-KL\alpha^*}$  is small (see Fig. 1b). This is because  $f_l$   
 51 will reach a local minima when  $e^{-KL\alpha} \rightarrow 0$ . If  $e^{-KL\alpha^*}$  is not large enough, it may be higher than  
 52  $\lim_{e^{-KL\alpha} \rightarrow 0} f_l(e^{-KL\alpha})$ , which means the shaded region near y-axis is also included. In this region  
 53  $f(\alpha)$  could be really small (which is the goal when optimizing the surrogate objective  $f(\alpha)$ ), but  
 54  $KL_{\alpha}$  could be extremely large.

55 To avoid this situation, we only have to assume that

$$f_u(e^{-KL\alpha^*}) \leq \lim_{e^{-KL\alpha} \rightarrow 0} f_l(e^{-KL\alpha}) = -\log \frac{C_4}{Z(\bar{S}; \alpha)}, \quad (\text{A13})$$

56 or if we denote  $\gamma_1 = \min_{\alpha} Z(\bar{S}; \alpha)$  and  $\gamma_2 = \max_{\alpha} Z(\bar{S}; \alpha)$ , then we only need the following  
 57 assumption:

58 **Assumption 2.** The optimal  $KL_{\alpha^*}$  is small so that  $f_u(e^{-KL\alpha^*}) \leq -\log \frac{C_4}{\gamma_1}$ .

59 This assumption holds as long as there is at least one task that is related to the main task (having  
 60 a small  $KL_{\alpha}$ ), which is reasonable because if all the tasks are unrelated, then reweighing is also  
 61 meaningless. See the remark below for more discussion on the validity of the assumption.

62 Now we give the formal version of the theorem:

63 **Theorem 1. (formal version)** With Assumption 1, 2, if  $\gamma_2 \leq \min(\frac{C_3}{e}, \frac{C_4}{e})$ , then we have

$$KL_{\hat{\alpha}} \leq KL_{\alpha^*} + \frac{2\gamma_2^2}{C} \log \frac{C'}{\gamma_2} \quad (\text{A14})$$

64 *Proof.* From Assumption 2 we have

$$\frac{C_2 e^{-KL_{\alpha^*}}}{C_2 e^{-KL_{\alpha^*}} + Z(\bar{S}; \alpha^*)} (KL_{\alpha^*} + C_1) + \frac{Z(\bar{S}; \alpha^*)}{C_2 e^{-KL_{\alpha^*}} + Z(\bar{S}; \alpha^*)} \log \frac{C_3}{Z(\bar{S}; \alpha^*)} \leq -\log \frac{C_4}{\gamma_1}. \quad (\text{A15})$$

65 Since  $\gamma_2 \leq C_3$  and  $\gamma_2 \leq C_4$ , we have  $\log \frac{C_3}{Z(\bar{S}; \alpha^*)} \geq \log \frac{C_3}{\gamma_2} \geq 0$ , and  $-\log \frac{C_4}{\gamma_1} \leq -\log \frac{C_4}{\gamma_2} \leq 0$ .  
 66 Then leaves us  $KL_{\alpha^*} + C_1 \leq 0$  in order to make (A15) satisfied. Then we can relax (A15) into

$$KL_{\alpha^*} + C_1 \leq -\log \frac{C_4}{\gamma_1}, \quad (\text{A16})$$

67 which gives

$$C_2 e^{-KL_{\alpha^*}} \geq \frac{C_5}{\gamma_1}, \quad (\text{A17})$$

68 where  $C_5 = C_4 C_2 e^{C_1}$ . This bounds the value of  $KL_{\alpha^*}$ .

69 Moreover, from (A12) and Assumption 2 we have

$$f_l(e^{-KL_{\hat{\alpha}}}) \leq f_u(e^{-KL_{\alpha^*}}) \leq -\log \frac{C_4}{\gamma_1}, \quad (\text{A18})$$

70 which gives

$$f_l(e^{-KL_{\hat{\alpha}}}) = \frac{C_2 e^{-KL_{\hat{\alpha}}}}{C_2 e^{-KL_{\hat{\alpha}}} + Z(\bar{S}; \hat{\alpha})} (KL_{\hat{\alpha}} + C_1) - \frac{Z(\bar{S}; \hat{\alpha})}{C_2 e^{-KL_{\hat{\alpha}}} + Z(\bar{S}; \hat{\alpha})} \log \frac{C_4}{Z(\bar{S}; \hat{\alpha})} \leq -\log \frac{C_4}{\gamma_1}. \quad (\text{A19})$$

71 Since  $Z(\bar{S}; \hat{\alpha}) \geq \gamma_1$ , we can relax (A19) into

$$\frac{C_2 e^{-KL_{\hat{\alpha}}}}{C_2 e^{-KL_{\hat{\alpha}}} + Z(\bar{S}; \hat{\alpha})} (KL_{\hat{\alpha}} + C_1) - \frac{Z(\bar{S}; \hat{\alpha})}{C_2 e^{-KL_{\hat{\alpha}}} + Z(\bar{S}; \hat{\alpha})} \log \frac{C_4}{Z(\bar{S}; \hat{\alpha})} \leq -\log \frac{C_4}{Z(\bar{S}; \hat{\alpha})}, \quad (\text{A20})$$

72 which can be simplified into

$$KL_{\hat{\alpha}} + C_1 \leq -\log \frac{C_4}{Z(\bar{S}; \hat{\alpha})} \leq -\log \frac{C_4}{\gamma_2}, \quad (\text{A21})$$

73 which means

$$C_2 e^{-KL_{\hat{\alpha}}} \geq \frac{C_5}{\gamma_2}. \quad (\text{A22})$$

74 This bounds the value of  $KL_{\hat{\alpha}}$ .

75 Now we build the connection between  $KL_{\hat{\alpha}}$  and  $KL_{\alpha^*}$ . Since  $f_l(e^{-KL_{\hat{\alpha}}}) \leq f_u(e^{-KL_{\alpha^*}})$ , we have

$$\begin{aligned} & \frac{C_2 e^{-KL_{\hat{\alpha}}}}{C_2 e^{-KL_{\hat{\alpha}}} + Z(\bar{S}; \hat{\alpha})} (KL_{\hat{\alpha}} + C_1) - \frac{Z(\bar{S}; \hat{\alpha})}{C_2 e^{-KL_{\hat{\alpha}}} + Z(\bar{S}; \hat{\alpha})} \log \frac{C_4}{Z(\bar{S}; \hat{\alpha})} \\ & \leq \frac{C_2 e^{-KL_{\alpha^*}}}{C_2 e^{-KL_{\alpha^*}} + Z(\bar{S}; \alpha^*)} (KL_{\alpha^*} + C_1) + \frac{Z(\bar{S}; \alpha^*)}{C_2 e^{-KL_{\alpha^*}} + Z(\bar{S}; \alpha^*)} \log \frac{C_3}{Z(\bar{S}; \alpha^*)}. \end{aligned} \quad (\text{A23})$$

76 Since  $KL_{\hat{\alpha}} + C_1 \leq -\log \frac{C_4}{\gamma_2} \leq 0$ ,  $KL_{\alpha^*} \geq 0$ , and also with (A17) and (A22), we can relax (A23)  
 77 into

$$\begin{aligned} & KL_{\hat{\alpha}} + C_1 - \frac{Z(\bar{S}; \hat{\alpha})}{C_5/\gamma_2} \log \frac{C_4}{Z(\bar{S}; \hat{\alpha})} \\ & \leq KL_{\alpha^*} + \frac{C_2 e^{-KL_{\alpha^*}}}{C_2 e^{-KL_{\alpha^*}} + Z(\bar{S}; \alpha^*)} C_1 + \frac{Z(\bar{S}; \alpha^*)}{C_5/\gamma_1} \log \frac{C_3}{Z(\bar{S}; \alpha^*)}, \end{aligned} \quad (\text{A24})$$

78 which gives

$$KL_{\hat{\alpha}} \leq KL_{\alpha^*} - \frac{Z(\bar{S}; \alpha^*)}{C_2 e^{-KL_{\alpha^*}} + Z(\bar{S}; \alpha^*)} C_1 + \frac{Z(\bar{S}; \hat{\alpha})}{C_5/\gamma_2} \log \frac{C_4}{Z(\bar{S}; \hat{\alpha})} + \frac{Z(\bar{S}; \alpha^*)}{C_5/\gamma_1} \log \frac{C_3}{Z(\bar{S}; \alpha^*)}. \quad (\text{A25})$$

79 Since  $Z(\bar{S}; \hat{\alpha}) \leq \gamma_2 \leq \frac{C_4}{e}$ , we have  $Z(\bar{S}; \hat{\alpha}) \log \frac{C_4}{Z(\bar{S}; \hat{\alpha})} \leq \gamma_2 \log \frac{C_4}{\gamma_2}$ . Similarly, we have  
 80  $Z(\bar{S}; \alpha^*) \log \frac{C_3}{Z(\bar{S}; \alpha^*)} \leq \gamma_2 \log \frac{C_3}{\gamma_2}$ . Then we have

$$KL_{\hat{\alpha}} \leq KL_{\alpha^*} - \frac{Z(\bar{S}; \alpha^*)}{C_2 e^{-KL_{\alpha^*}} + Z(\bar{S}; \alpha^*)} C_1 + \frac{\gamma_2^2}{C_5} \log \frac{C_4}{\gamma_2} + \frac{\gamma_2^2}{C_5} \log \frac{C_3}{\gamma_2}. \quad (\text{A26})$$

81 Since  $C_1 \leq 0$ , we can get

$$KL_{\hat{\alpha}} \leq KL_{\alpha^*} + \frac{\gamma_2^2}{C_5} (-C_1) + \frac{\gamma_2^2}{C_5} \log \frac{C_4}{\gamma_2} + \frac{\gamma_2^2}{C_5} \log \frac{C_3}{\gamma_2}, \quad (\text{A27})$$

82 which gives

$$KL_{\hat{\alpha}} \leq KL_{\alpha^*} + \frac{2\gamma_2^2}{C_5} \log \frac{C_6}{\gamma_2}, \quad (\text{A28})$$

83 where  $C_6 = \sqrt{C_3 C_4 e^{-C_1}}$ .

84

□

85 **Remark.** From Theorem 1 we see that  $KL_{\hat{\alpha}}$  is close to  $KL_{\alpha^*}$  as long as  $\gamma_2$  is small. One may  
86 notice that  $\gamma_2$  cannot be arbitrarily small because from (A22) we have

$$\frac{C_5}{\gamma_2} \leq C_2 e^{-KL_{\hat{\alpha}}} \leq C_2, \quad (\text{A29})$$

87 which means

$$\gamma_2 \geq \frac{C_5}{C_2} = C_4 e^{C_1}. \quad (\text{A30})$$

88 However, we can safely assume that

$$C_1 = \log \frac{p^m(\theta^*)}{p^*(\theta^*)} \ll 0 \quad (\text{A31})$$

89 since  $p^*$  is much more informative than  $p^m$ , especially when labeled data for the main task is scarce.  
90 This means  $\gamma_2$  can be extremely small as long as  $C_1$  is small, which makes  $KL_{\hat{\alpha}}$  close to  $KL_{\alpha^*}$ .  
91 Similarly, Assumption 2 can easily hold as long as  $C_1$  is small.

## 92 1.2 Sampling through Langevin Dynamics (P2)

93 In **Samples (P2)** we use Langevin dynamics [16, 22] to sample from the distribution  $p^J$ . Concretely,  
94 at each iteration, we update  $\theta$  by

$$\theta_{t+1} = \theta_t - \epsilon_t \nabla \mathcal{L}(\theta_t) + \eta_t, \quad (\text{A32})$$

95 where  $\mathcal{L}(\theta) \propto -\log p^J(\theta)$  is the joint loss, and  $\eta_t \sim N(0, 2\epsilon_t)$  is a Gaussian noise. In this way,  $\theta_t$   
96 converges to samples from  $p^J$ , which can be used to estimate our optimization objective. However,  
97 since we normally use a mini-batch estimator  $\hat{\mathcal{L}}(\theta)$  to approximate  $\mathcal{L}(\theta)$ , this may introduce additional  
98 noise other than  $\eta_t$ , which may make the sampling procedure inaccurate. In [22] it is proposed to  
99 anneal the learning rate to zero so that the gradient stochasticity is dominated by the injected noise,  
100 thus alleviating the impact of mini-batch estimator. However we find in practice that the gradient  
101 noise is negligible compared to the injected noise (Table 1). Therefore, we ignore the gradient noise  
102 and directly inject the noise  $\eta_t$  into the updating step.

Table 1: Standard deviation of different types of noise. We find that the gradient noise is negligible compared to the injected noise.

Standard deviation	
Gradient Noise	$\sim 10^{-6}$
Injected Noise	$\sim 10^{-3}$

## 103 1.3 Score Function and Fisher Divergence (P3)

104 In **Partition Function (P3)** we propose to minimize

$$\min_{\alpha} E_{\theta \sim p^J} \|\nabla \log p(\mathcal{T}_m | \theta) - \nabla \log p_{\alpha}(\theta)\|_2^2 \quad (\text{A33})$$

105 as our final objective. Notice that

$$\begin{aligned} & \min_{\alpha} E_{\theta \sim p^J} \|\nabla \log p(\mathcal{T}_m | \theta) - \nabla \log p_{\alpha}(\theta)\|_2^2 \\ \Leftrightarrow & \min_{\alpha} E_{\theta \sim p^J} \|\nabla \log p^m(\theta) - \nabla \log p_{\alpha}(\theta)\|_2^2 \\ \Leftrightarrow & \min_{\alpha} E_{\theta \sim p^J} \|\nabla \log(p^m(\theta) \cdot p_{\alpha}(\theta)) - 2 \cdot \nabla \log p_{\alpha}(\theta)\|_2^2 \\ \Leftrightarrow & \min_{\alpha} E_{\theta \sim p^J} \|\nabla \log p^J(\theta) - \nabla \log p_{\alpha}^2(\theta)\|_2^2 \\ \Leftrightarrow & \min_{\alpha} F(p^J(\theta) \parallel \frac{1}{Z'(\alpha)} p_{\alpha}^2(\theta)), \end{aligned} \quad (\text{A34})$$

106 where  $F(p(\theta) \parallel q(\theta)) = E_{\theta \sim p} \|\nabla \log p(\theta) - \nabla \log q(\theta)\|_2^2$  is the *Fisher divergence*, and  $Z'(\boldsymbol{\alpha}) =$   
 107  $\int p_{\boldsymbol{\alpha}}^2(\theta) d\theta$  is the normalization term. This means, by optimizing (A33), we are actually minimizing  
 108 the Fisher divergence between  $p^J(\theta)$  and  $\frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta)$ . As pointed by [8, 13], Fisher divergence  
 109 is stronger than KL divergence, which means by minimizing  $F(p^J(\theta) \parallel \frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta))$ , the KL  
 110 divergence  $D_{KL}(p^J(\theta) \parallel \frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta))$  is also bounded near the optimum up to a small error.

111 Therefore, optimizing (A33) is equivalent to minimizing  $D_{KL}(p^J(\theta) \parallel \frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta))$ . Notice that

$$\begin{aligned}
 & \min_{\boldsymbol{\alpha}} D_{KL}(p^J(\theta) \parallel \frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta)) \\
 \Leftrightarrow & \min_{\boldsymbol{\alpha}} \int p^J(\theta) \log \frac{p^J(\theta)}{\frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta)} d\theta \\
 \Leftrightarrow & \min_{\boldsymbol{\alpha}} \int p^J(\theta) \log \frac{\frac{1}{Z(\boldsymbol{\alpha})} p^m(\theta) p_{\boldsymbol{\alpha}}(\theta)}{\frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta)} d\theta \tag{A35} \\
 \Leftrightarrow & \min_{\boldsymbol{\alpha}} \int p^J(\theta) \log \frac{p^m(\theta)}{p_{\boldsymbol{\alpha}}(\theta)} d\theta + \log \frac{Z'(\boldsymbol{\alpha})}{Z(\boldsymbol{\alpha})} \\
 \Leftrightarrow & \min_{\boldsymbol{\alpha}} \int p^J(\theta) \log \frac{p^m(\theta)}{p_{\boldsymbol{\alpha}}(\theta)} d\theta + \log \frac{\int p_{\boldsymbol{\alpha}}^2(\theta) d\theta}{\int p^m(\theta) p_{\boldsymbol{\alpha}}(\theta) d\theta}
 \end{aligned}$$

112 is different from (A1) only on the  $\log \frac{\int p_{\boldsymbol{\alpha}}^2(\theta) d\theta}{\int p^m(\theta) p_{\boldsymbol{\alpha}}(\theta) d\theta}$  term. To analyze the impact of this additional  
 113 term, we assume that the likelihood function of each auxiliary task is a Gaussian, *i.e.*,  $p(\mathcal{T}_{\alpha_k} | \theta) \propto$   
 114  $N(\theta | \theta_k, \boldsymbol{\Sigma})$ , with mean  $\theta_k$  and covariance  $\boldsymbol{\Sigma}$ . Then we have  $p_{\boldsymbol{\alpha}}(\theta) = N(\theta | \sum_k \alpha_k \theta_k / K, \boldsymbol{\Sigma} / K)$   
 115 (note that  $\sum_k \alpha_k = K$ ). In this case  $\int p_{\boldsymbol{\alpha}}^2(\theta) d\theta$  only depends on  $\boldsymbol{\Sigma}$  and is invariant to  $\boldsymbol{\alpha}$ . Thus  
 116 optimizing (A33) is equivalent to

$$\begin{aligned}
 & \min_{\boldsymbol{\alpha}} D_{KL}(p^J(\theta) \parallel \frac{1}{Z'(\boldsymbol{\alpha})} p_{\boldsymbol{\alpha}}^2(\theta)) \\
 \Leftrightarrow & \min_{\boldsymbol{\alpha}} \int p^J(\theta) \log \frac{p^m(\theta)}{p_{\boldsymbol{\alpha}}(\theta)} d\theta + \log \frac{\int p_{\boldsymbol{\alpha}}^2(\theta) d\theta}{\int p^m(\theta) p_{\boldsymbol{\alpha}}(\theta) d\theta} \tag{A36} \\
 \Leftrightarrow & \min_{\boldsymbol{\alpha}} \int p^J(\theta) \log \frac{p^m(\theta)}{p_{\boldsymbol{\alpha}}(\theta)} d\theta - \log \int p^m(\theta) p_{\boldsymbol{\alpha}}(\theta) d\theta.
 \end{aligned}$$

117 Denote the optimal solution for (A36) by  $\boldsymbol{\alpha}^\dagger$ . Then we can build the connection between  $\boldsymbol{\alpha}^\dagger$  and  $\hat{\boldsymbol{\alpha}}$   
 118 by

$$\int p^J(\theta) \log \frac{p^m(\theta)}{p_{\boldsymbol{\alpha}^\dagger}(\theta)} d\theta - \log \int p^m(\theta) p_{\boldsymbol{\alpha}^\dagger}(\theta) d\theta \leq \int p^J(\theta) \log \frac{p^m(\theta)}{p_{\hat{\boldsymbol{\alpha}}}(\theta)} d\theta - \log \int p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta. \tag{A37}$$

119 Since  $\hat{\boldsymbol{\alpha}}$  minimizes  $\int p^J(\theta) \log \frac{p^m(\theta)}{p_{\boldsymbol{\alpha}}(\theta)} d\theta$ , which means  $\int p^J(\theta) \log \frac{p^m(\theta)}{p_{\hat{\boldsymbol{\alpha}}}(\theta)} d\theta \leq \int p^J(\theta) \log \frac{p^m(\theta)}{p_{\boldsymbol{\alpha}^\dagger}(\theta)} d\theta$ ,  
 120 we can get

$$-\log \int p^m(\theta) p_{\boldsymbol{\alpha}^\dagger}(\theta) d\theta \leq -\log \int p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta, \tag{A38}$$

121 or

$$\int p^m(\theta) p_{\boldsymbol{\alpha}^\dagger}(\theta) d\theta \geq \int p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta, \tag{A39}$$

122 which gives

$$\int_{\theta \in S} p^m(\theta) p_{\boldsymbol{\alpha}^\dagger}(\theta) d\theta + \int_{\theta \in \bar{S}} p^m(\theta) p_{\boldsymbol{\alpha}^\dagger}(\theta) d\theta \geq \int_{\theta \in S} p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta + \int_{\theta \in \bar{S}} p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta. \tag{A40}$$

123 Then we have

$$\begin{aligned}
 \int_{\theta \in S} p^m(\theta) p_{\boldsymbol{\alpha}^\dagger}(\theta) d\theta & \geq \int_{\theta \in S} p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta + \int_{\theta \in \bar{S}} p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta - \int_{\theta \in \bar{S}} p^m(\theta) p_{\boldsymbol{\alpha}^\dagger}(\theta) d\theta \\
 & \geq \int_{\theta \in S} p^m(\theta) p_{\hat{\boldsymbol{\alpha}}}(\theta) d\theta - (\gamma_2 - \gamma_1).
 \end{aligned} \tag{A41}$$

124 From Assumption 1 we have

$$\frac{p^m(\theta^*) p_{\boldsymbol{\alpha}^\dagger}(\theta^*)}{p^*(\theta^*)} \geq \frac{p^m(\theta^*) p_{\hat{\boldsymbol{\alpha}}}(\theta^*)}{p^*(\theta^*)} - (\gamma_2 - \gamma_1), \tag{A42}$$

125 which gives

$$KL_{\alpha^\dagger} = -\log \frac{p_{\alpha^\dagger}(\theta^*)}{p^*(\theta^*)} \leq -\log \left( \frac{p_{\hat{\alpha}}(\theta^*)}{p^*(\theta^*)} - \frac{\gamma_2 - \gamma_1}{p^m(\theta^*)} \right) \leq -\log \frac{p_{\hat{\alpha}}(\theta^*)}{p^*(\theta^*)} + \frac{\gamma_2 - \gamma_1}{p^m(\theta^*)}, \quad (\text{A43})$$

126 or

$$KL_{\alpha^\dagger} \leq KL_{\hat{\alpha}} + \frac{\gamma_2}{C_2}. \quad (\text{A44})$$

127 After combining with Theorem 1, we have

$$KL_{\alpha^\dagger} \leq KL_{\alpha^*} + \frac{2\gamma_2^2}{C_5} \log \frac{C_6}{\gamma_2} + \frac{\gamma_2}{C_2}. \quad (\text{A45})$$

128 This means by optimizing our final objective (A33), the KL divergence  $KL_{\alpha^\dagger}$  is also bounded near  
129 the optimal value, which provides a theoretical justification of our algorithm.

## 130 1.4 Tips for Practitioners

131 In Section 2.4, we propose a two-stage algorithm, where we update the task weights with Langevin  
132 dynamics in the first stage, and then update the model with fixed task weights in the second stage.  
133 However, we find in practice that we can also find the similar task weights if we turn off the Langevin  
134 dynamics and directly sample from regular SGD. Therefore, we can further simplify the algorithm  
135 by removing the Langevin dynamics and merge the two stage, *i.e.*, update task weights and model  
136 parameters at the same time until convergence. This simplified version is summarized in Algorithm 1.

---

### Algorithm 1 ARML (simplified version)

---

**Input:** main task data  $\mathcal{T}_m$ , auxiliary task data  $\mathcal{T}_{a_k}$ , initial parameter  $\theta_0$ , initial task weights  $\alpha$   
**Parameters:** learning rate of  $t$ -th iteration  $\epsilon_t$ , learning rate for task weights  $\beta$

**for** iteration  $t = 1$  to  $T$  **do**

$$\theta_t \leftarrow \theta_{t-1} - \epsilon_t (-\nabla \log p(\mathcal{T}_m | \theta_{t-1}) - \sum_{k=1}^K \alpha_k \nabla \log p(\mathcal{T}_{a_k} | \theta_{t-1})) + \eta_t$$

$$\alpha \leftarrow \alpha - \beta \nabla_{\alpha} \|\nabla \log p(\mathcal{T}_m | \theta_t) - \sum_{k=1}^K \alpha_k \nabla \log p(\mathcal{T}_{a_k} | \theta_t)\|_2^2$$

Project  $\alpha$  back into  $\mathcal{A}$

**end for**

---

## 137 2 Experimental Settings

138 For all results, we repeat experiments for three times and report the average performance. Error bars  
139 are reported with CI=95%. In our algorithm, the only hyperparameter is the learning rate  $\beta$  of task  
140 weights. Specifically, we find the results insensitive to the choice of  $\beta$ . Therefore, we randomly  
141 choose  $\beta \in [0.0005, 0.05]$ , for a trade-off between steady training and fast convergence. We use  
142 PyTorch [19] for implementation.

### 143 2.1 Semi-supervised Learning

144 For semi-supervised learning, we use two datasets, CIFAR10 [11] and SVHN [17]. For CIFAR10,  
145 we follow the standard train/validation split, with 45000 images for training and 5000 for validation.  
146 Only 4000 out of 45000 training images are labeled. For SVHN, we use the standard train/validation  
147 split with 65932 images for training and 7325 for validation. Only 1000 out of 65392 images are  
148 labeled. Both datasets can be downloaded from the official PyTorch torchvision library (<https://pytorch.org/docs/stable/torchvision/index.html>). Following [18], we use WRN-28-2  
149 as our backbone, *i.e.*, ResNet [7] with depth 28 and width 2, including batch normalization [9] and  
150 leaky ReLU [15]. We train our model for 200000 iterations, using Adam [10] optimizer with batch  
151 size of 256 and learning rate of 0.005 in first 160000 iterations and 0.001 for the rest iterations.  
152

153 For implementation of self-supervised semi-supervised learning (S4L), we follow the settings in  
154 the original paper [23]. Note that we make two differences from [23]: (i) for steadier training, we  
155 use the model with time-averaged parameters [21] to extract feature of the original image, (ii) To  
156 avoid over-sampling of negative samples in triplet-loss [1], we only put a loss on the cosine similarity  
157 between original feature and augmented feature.



## 158 2.2 Multi-label Classification

159 For multi-label classification, we use CelebA [14] as our dataset. It contains 200K face images,  
160 each labeled with 40 binary attributes. We cast this into a multi-label classification problem, where  
161 we randomly choose one attribute as the main classification task, and other 39 as auxiliary tasks.  
162 We randomly choose 1% images as labeled images for main task. The dataset is available at <http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>. We use ResNet18 [7] as our backbone. We  
164 train the model for 90 epochs using SGD solver with batch size of 256 and scheduled learning rate of  
165 0.1 initially and  $0.1 \times$  shrunked every 30 epochs.

## 166 2.3 Domain Generalization

167 Following the literature [2, 4], we use PACS [12] as our dataset for domain generalization. PACS  
168 consists of four domains (photo, art painting, cartoon and sketch), each containing 7 categories (dog,  
169 elephant, giraffe, guitar, horse, house and person). The dataset is created by intersecting classes  
170 in Caltech-256 [6], Sketchy [20], TU-Berlin [5] and Google Images. Dataset can be downloaded  
171 from <http://sketchx.eecs.qmul.ac.uk/>. Following protocol in [12], we split the images from  
172 training domains to 9 (train) : 1 (val) and test on the whole target domain. We use a simple data  
173 augmentation protocol by randomly cropping the images to 80-100% of original sizes and randomly  
174 apply horizontal flipping. We use ResNet18 [7] as our backbone. Models are trained with SGD  
175 solver, 100 epochs, batch size 128. Learning rate is set to 0.001 and shrunked down to 0.0001 after 80  
176 epochs.

## 177 References

- 178 [1] Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A  
179 theoretical analysis of contrastive unsupervised representation learning. *arXiv preprint arXiv:1902.09229*,  
180 2019.
- 181 [2] Nader Asadi, Mehrdad Hosseinzadeh, and Mahdi Eftekhari. Towards shape biased unsupervised represen-  
182 tation learning for domain generalization. *arXiv preprint arXiv:1909.08245*, 2019.
- 183 [3] Jonathan Baxter. A bayesian/information theoretic model of learning to learn via multiple task sampling.  
184 *Machine learning*, 28(1):7–39, 1997.
- 185 [4] Fabio M Carlucci, Antonio D’Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. Domain  
186 generalization by solving jigsaw puzzles. In *Proceedings of the IEEE Conference on Computer Vision and*  
187 *Pattern Recognition*, pages 2229–2238, 2019.
- 188 [5] Mathias Eitz, James Hays, and Marc Alexa. How do humans sketch objects? *ACM Transactions on*  
189 *graphics (TOG)*, 31(4):1–10, 2012.
- 190 [6] Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007.
- 191 [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
192 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 193 [8] Tianyang Hu, Zixiang Chen, Hanxi Sun, Jincheng Bai, Mao Ye, and Guang Cheng. Stein neural sampler.  
194 *arXiv preprint arXiv:1810.03545*, 2018.
- 195 [9] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing  
196 internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- 197 [10] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
198 *arXiv:1412.6980*, 2014.
- 199 [11] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 200 [12] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain  
201 generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550,  
202 2017.
- 203 [13] Qiang Liu, Jason Lee, and Michael Jordan. A kernelized stein discrepancy for goodness-of-fit tests. In  
204 *International conference on machine learning*, pages 276–284, 2016.



- 205 [14] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In  
206 *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015.
- 207 [15] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network  
208 acoustic models. In *Proc. icml*, volume 30, page 3, 2013.
- 209 [16] Radford M Neal et al. Mcmc using hamiltonian dynamics. *Handbook of markov chain monte carlo*,  
210 2(11):2, 2011.
- 211 [17] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits  
212 in natural images with unsupervised feature learning. 2011.
- 213 [18] Avital Oliver, Augustus Odena, Colin A Raffel, Ekin Dogus Cubuk, and Ian Goodfellow. Realistic  
214 evaluation of deep semi-supervised learning algorithms. In *Advances in Neural Information Processing*  
215 *Systems*, pages 3235–3246, 2018.
- 216 [19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,  
217 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep  
218 learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.
- 219 [20] Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. The sketchy database: learning to retrieve  
220 badly drawn bunnies. *ACM Transactions on Graphics (TOG)*, 35(4):1–12, 2016.
- 221 [21] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency  
222 targets improve semi-supervised deep learning results. In *Advances in neural information processing*  
223 *systems*, pages 1195–1204, 2017.
- 224 [22] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings*  
225 *of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011.
- 226 [23] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4l: Self-supervised semi-supervised  
227 learning. In *Proceedings of the IEEE international conference on computer vision*, pages 1476–1485,  
228 2019.