

A. Experiments: learning hyperparameters

We train task model with the following hyperparameters for all methods and datasets: SGD optimizer with 0.9 Nesterov momentum, $1e-4$ weight decay, 200 train epochs, 256 mini-batch size, cosine learning rate annealing with 5 warm-up epochs. Learning rate is 0.005 for SVHN and 0.1 for CIFAR-10/100 and ImageNet. The HO optimizer is RMSprop with 0.01 learning rate for ImageNet and 0.05 for others, and it uses the same schedule as the task optimizer.

Table 1. Runtime comparison of our and recently published methods, GPU hours. AutoDO results are estimated using V100 GPU.

Alg./Dataset	CIFAR-10	SVHN
AA [3]	5,000	1,000
PBA [5]	5.0	1.0
FAA [8]	3.5	1.5
OHL [9]	83.3	-
DADA [7]	0.1	0.1
AutoDO, 1 epoch	0.036	0.046
AutoDO, $E = 150$	1.8	2.3
AutoDO, $E = 50$	5.4	6.8

We use only θ of penultimate fully-connected layer in (5) with $T = 5$ Neumann series iterations for \mathbf{H}_θ^{-1} approximation. The T hyperparameter is chosen from Fig. 3 ablation study in [10]. The hyperparameter E in Alg. 1 is 100 for ImageNet and 50 for others. The soft-label initialization constant from Section 4.2 is $\alpha = 0.1$. The evaluated architectures are ResNet-18 for ImageNet, Wide-ResNet-28-10 for CIFAR-10/100 and SVHN. We replace ReLU nonlinearities in all networks with CELU [1] to satisfy C^1 requirement in (4). While IFT in (4) is only locally defined ($\|\lambda - \hat{\lambda}\| \leq r_1$ and $\|\theta - \hat{\theta}\| \leq r_2$) for a fixed point $(\hat{\lambda}, \hat{\theta})$, practically, the choice of hyperparameter E and CELU stabilizes optimization process in Alg. 1. In addition, there are an attempts to extend IFT to a global case in [6, 2, 4].

We run all experiments on P100/V100 GPUs and each one takes few hours depending on the dataset setting. Only large-scale ImageNet experiments can take up to several days. Table 1 presents detailed comparison between the estimated AutoDO runtime and the reported runtimes for other recent methods.

B. Experiments: SVHN learning curves

The additional learning curves for SVHN are showed in Figure 1. Our AutoDO optimization starts at epoch $E = \{10, 50, 150, 190\}$. It is evident that $E = 50$ hyperparameter has a minor edge over other settings, but practically $E = 150$ works almost the same with $3\times$ less computing.

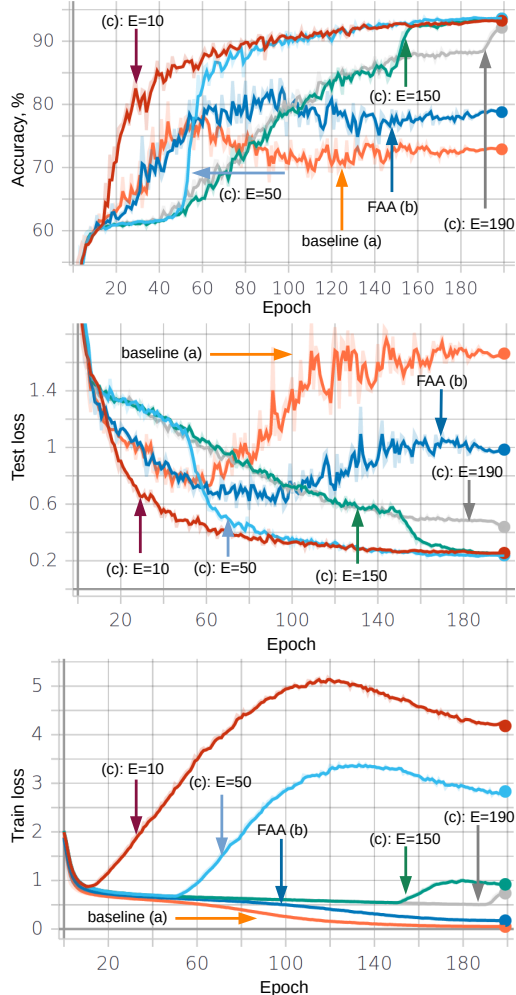


Figure 1. SVHN top-1 test accuracy (top), test loss (middle) and train loss (bottom) learning curves for the following models trained on a dataset with $100\times$ class imbalance and 10% label noise: (a) baseline, (b) FAA [8], and (c) our AutoDO with $\lambda^{A,W,S}$. Our HO starts at epoch $E = \{10, 50, 150, 190\}$, which prevents overfitting to the distorted train data and improves generalization to test data.

C. Theoretical background: implicit differentiation meets density matching

The gradient of expectation in (7) over two independent distributions \hat{Q}_x^{val} and \hat{Q}_x is a product

$$\nabla_{\lambda} \mathbb{E}_{\hat{Q}_x^{\text{val}}, \hat{Q}_x} [\mathcal{L}_v] = -\mathbb{E}_{\hat{Q}_x^{\text{val}}} [\nabla_{\theta} \mathcal{L}_v] \mathbb{E}_{\hat{Q}_x} [\mathbf{H}_\theta^{-1} \nabla_{\theta} \nabla_{\lambda}^T \mathcal{L}].$$

The Hessian \mathbf{H}_θ itself is calculated as expectation over \hat{Q}_x due to computational and stability reasons. It leads to a connection to the Fisher information metric \mathcal{I}_θ [11] as

$$\mathbb{E}_{\hat{Q}_x} [\mathbf{H}_\theta] = \mathbb{E}_{\hat{Q}_x} \left[\frac{\partial^2 \mathcal{L}}{\partial \theta \partial \theta^T} \right] = \mathbb{E}_{\hat{Q}_x} [\mathbf{u}_i(\theta) \mathbf{u}_i(\theta)^T] = -\mathcal{I}_\theta,$$

where $\mathbf{u}_i(\theta) = -\partial \mathcal{L}(i) / \partial \theta = \nabla_{\theta} \log p(\mathbf{y}_i | \mathbf{x}_i, \theta(\lambda))$ are

the Fisher scores [12] for log-likelihood loss function.

The expectation $\mathbb{E}_{\hat{Q}_x} [\nabla_{\theta} \nabla_{\lambda}^T \mathcal{L}]$ contains a second-order derivative of the train loss $\mathcal{L}(i) = -\log p(\mathbf{y}_i | \mathbf{x}_i, \theta(\lambda)) = -\log p$ from (6) that can be similarly simplified to

$$\begin{aligned} \mathbb{E}_{\hat{Q}_x} \left[\frac{\partial^2 \mathcal{L}}{\partial \theta \partial \lambda^T} \right] &= -\mathbb{E}_{\hat{Q}_x} \left[\frac{\partial}{\partial \theta} \left(\frac{\partial \log p}{\partial \lambda^T} \right) \right] = \\ &= -\mathbb{E}_{\hat{Q}_x} \left[-\frac{1}{p^2} \frac{\partial p}{\partial \theta} \frac{\partial p}{\partial \lambda^T} + \frac{1}{p} \frac{\partial^2 p}{\partial \theta \partial \lambda^T} \right] = \mathbb{E}_{\hat{Q}_x} \left[\frac{\partial \mathcal{L}}{\partial \theta} \frac{\partial \mathcal{L}}{\partial \lambda^T} \right], \end{aligned}$$

where the term $-\mathbb{E}_{\hat{Q}_x} \left[\frac{1}{p} \frac{\partial^2 p}{\partial \theta \partial \lambda^T} \right] = -\frac{\partial^2}{\partial \theta \partial \lambda^T} \int_x p d\mathbf{x} = 0$.

By substituting the above derivations, the final form of (7) and its practical variant with equal-probability data points can be obtained as

$$\begin{aligned} \nabla_{\lambda} \mathbb{E}_{\hat{Q}_x^{\text{val}}, \hat{Q}_x} [\mathcal{L}_v] &= \mathbb{E}_{\hat{Q}_x^{\text{val}}} [\mathbf{u}^v(\theta)] \mathcal{I}_{\theta}^{-1} \mathbb{E}_{\hat{Q}_x} [\mathbf{u}(\theta) \mathbf{u}(\lambda)^T] \\ &= \left[\frac{1}{M} \sum_{j \in \mathbb{M}} \mathbf{u}_j^v(\theta) \right] \mathcal{I}_{\theta}^{-1} \left[\frac{1}{N} \sum_{i \in \mathbb{N}} \mathbf{u}_i(\theta) \mathbf{u}_i^T(\lambda) \right]. \end{aligned}$$

References

- [1] Jonathan T. Barron. Continuously differentiable exponential linear units. *arXiv:1704.07483*, 2017. 1
- [2] Mihai Cristea. On global implicit function theorem. *Journal of Mathematical Analysis and Applications*, 2017. 1
- [3] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. AutoAugment: Learning augmentation policies from data. *arXiv:1805.09501*, 2018. 1
- [4] M. Galewski and M. Rădulescu. On a global implicit function theorem for locally Lipschitz maps via nonsmooth critical point theory. *arXiv:1704.04280*, 2017. 1
- [5] Daniel Ho, Eric Liang, Ion Stoica, Pieter Abbeel, and Xi Chen. Population based augmentation: Efficient learning of augmentation policy schedules. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019. 1
- [6] Dariusz Idczak. A global implicit function theorem and its applications to functional equations. *Discrete and Continuous Dynamical Systems - Series B*, 2014. 1
- [7] Yonggang Li, Guosheng Hu, Yongtao Wang, Timothy Hospedales, Neil M. Robertson, and Yongxin Yang. DADA: Differentiable automatic data augmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 1
- [8] Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim. Fast AutoAugment. In *Advances in Neural Information Processing Systems*, 2019. 1
- [9] Chen Lin, Minghao Guo, Chuming Li, Xin Yuan, Wei Wu, Junjie Yan, Dahua Lin, and Wanli Ouyang. Online hyperparameter learning for auto-augmentation strategy. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 1
- [10] Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters by implicit differentiation. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020. 1
- [11] Alexander Ly, Maarten Marsman, A J Verhagen, Raoul Grasman, and Eric-Jan Wagenmakers. A tutorial on Fisher information. *Journal of Mathematical Psychology*, 2017. 1
- [12] Laurens van der Maaten. Learning discriminative Fisher kernels. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2011. 2