An Online Method for A Class of Distributionally Robust Optimization with Non-Convex Objectives

Qi Qi^{*†}, Zhishuai Guo^{*†}, Yi Xu[‡], Rong Jin[‡], Tianbao Yang[†] [†] Department of Computer Science, University of Iowa, Iowa City, IA 52242 [‡] Machine Intelligence Technology, Alibaba Group {qi-qi, zhishuai-guo, tianbao-yang}@uiowa.edu,{yixu, jinrong.jr}@alibaba-inc.com * Equal contribution

Background and Problem Definition

Distributionally Robust Optimization (DRO):

• Problem Formulation:

$$\min_{\mathbf{w}\in \mathbb{R}^d} \max_{\mathbf{p}\in\Delta_n} \left\{ F_p(\mathbf{w}) = \sum_{i=1}^n p_i \ell(\mathbf{w}; \mathbf{z}_i) - h(\mathbf{p}, \mathbf{1}/n) + r(\mathbf{w}; \mathbf{z}_i) - h(\mathbf{p}, \mathbf{1}/n) + r(\mathbf{w$$

• where $\Delta_n = \{ \mathbf{p} \in R^n : \sum_i p_i = 1, p_i \ge 0 \}.$ $\ell(\mathbf{w}; \mathbf{z})$: denotes a loss function on data $\mathbf{z} = (\mathbf{x}, y_I) \sim \mathbf{D}$, $h(\mathbf{p}, \mathbf{1}/n)$: a divergence measure $r(\mathbf{w})$ is convex regularizer of \mathbf{w}

Non-Convex Concave Min-Max Optimization

Related Work

Primal-Dual algorithms:

• PG-SMD [1], Epoch Primal-Dual SGD [2], $O(1/\epsilon^4)$

- Stoc-AGDA [3], PE-SGDA [4]
- $O(1/\mu^2 \varepsilon)$ by leveraging PL Condition

Deficiencies:

• O(n) computational cost per **p** variable updates.

Note: As the computational cost is intolerable for updating **p** variable is huge, Primal-Dual algorithms are not as favorable as SGD-type algo

Research Question

Can we design an efficient algorithm that is independent of data convergence rates and is also applicable to deep neural network trai

Compositional Equivalence of $F_{\mathbf{p}}(\mathbf{w})$

Considering the *KL*-divergence measure, e.g., $h(\mathbf{p}, \mathbf{1}/n) = -\sum p_i \log(n)$

$$\min_{\mathbf{w}\in\mathbb{R}^d} \left\{ F_{dro}(\mathbf{w}) = \lambda \log\left(\mathbb{E}_{\mathbf{z}} \exp\left(\frac{\ell(\mathbf{w};\mathbf{z})}{\lambda}\right)\right) + r(\mathbf{w})$$

General Compositional Optimization

Assumptions

Assumption 1: (Non-convex Setting) Let C_f, L_f, C_g and L_g be positive that

(a) $f: \mathbb{R}^p \to \mathbb{R}$ is a C_f -Lipschitz function and its gradient ∇f is L_f -Lips (b) $g_{\mathbf{z}} : \mathbb{R}^d \to \mathbb{R}^p$ satisfies $\mathbb{E} \|g_{\mathbf{z}}(\mathbf{w}_1) - g_{\mathbf{z}}(\mathbf{w}_2)\|^2 \le C_g^2 \|\mathbf{w}_1 - \mathbf{w}_2\|^2$ for an Jacobian $\nabla g_{\mathbf{z}}$ satisfies $\mathbb{E}[\|\nabla g_{\mathbf{z}}(\mathbf{w}_1) - \nabla g_{\mathbf{z}}(\mathbf{w}_2)\|^2] \leq L_g^2 \|\mathbf{w}_1 - \mathbf{w}_2\|^2$.

(c) $r: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ is a convex and lower-semicontinuous function.

(d) $F_* = \inf_{\mathbf{w}} F(\mathbf{w}) \ge -\infty$ and $F(\mathbf{w}_1) - F_* \le \Delta_F$ for the initial solution \mathbf{w}_1

(e) Let σ_g and $\sigma_{g'}$ be positive constants and $\sigma^2 = \sigma_g^2 + \sigma_{g'}^2$. Assume that

 $\mathbb{E}_{\mathbf{z}}[\|g_{\mathbf{z}}(\mathbf{w}) - g(\mathbf{w})\|^2] \le \sigma_g^2, \ \mathbb{E}_{\mathbf{z}}[\|\nabla g_{\mathbf{z}}(\mathbf{w}) - \nabla g(\mathbf{w})\|^2]$

Assumption 2: (PL Condition Setting) When $r(\mathbf{w})$ is a smooth function, $F(\mathbf{w})$ satisfies the μ -PL condition if there exists $\mu > 0$ such that

$$2\mu(F(\mathbf{w}) - \min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w})) \le \|\nabla F(\mathbf{w})\|^2.$$

Literature

		Proposed Stochastic Algorithm
		Proximal Gradient Convergence Measure:
		$\mathscr{G}_{\eta}(\mathbf{w}) = \frac{1}{n} (\mathbf{w} - \operatorname{prox}_{r}^{\eta} (\mathbf{w} - \mathbf{w}))$
$\mathbf{v})\Big\},$	(1)	When $r = 0$, the proximal gradient reduces to the $\nabla F(\mathbf{w})$.
		Sample Complexity: The sample complexity order to achieve $\mathbb{E}[\ \mathscr{G}_{\eta}(\mathbf{w})\ ^2] \leq \varepsilon$ for a certain 1
		Algorithm 1 COVER $(\mathbf{w}_1, u_1, \mathbf{v}_1, \{\eta_t\}, T, PL =$
		1: Let $a_t = c \eta_t^2$ 2: if not PL then
		 3: Draw a samples z and construct the estimat 4: end if
		5: for $t = 1,, T - 1$ do 6: $\mathbf{w}_{t+1} \leftarrow \operatorname{prox}_{r}^{\eta_{t}}(\mathbf{w}_{t} - \eta_{t}\mathbf{v}_{t}^{\top}\nabla f(u_{t}))$ 7: Draw a samples \mathbf{z}_{t+1} , and update
		$u_{t+1} = g_{\mathbf{z}_{t+1}}(\mathbf{w}_{t+1}) + (1$
		$\mathbf{v}_{t+1} = \nabla g_{\mathbf{z}_{t+1}}(\mathbf{w}_{t+1}) + (1$
when the d orithms for	lata size <i>n</i> DL.	8: end for 9: Return: $(\mathbf{w}_{\tau}, u_{\tau}, \mathbf{v}_{\tau})$ for randomly selected τ
		Algorithm 2 RECOVER $(\mathbf{w}_0, \boldsymbol{\varepsilon}_0, c)$
size n, ha	s faster	 Initialization: Draw a sample z₀ and con ∇g_{z₀}(w₀) for k = 1,,K do (w_k, u_k, v_k) = COVER(w_{k-1}, u_{k-1}, v_{k-1}, η_k, T_k change η_k, T_k according to Theorem RECOV end for Return: w_K
$g(np_i)$		Theorem COVER
$\Big)\Big\}.$	(2)	Assume Assumption 1 holds, for any $C > 0$,
		$\max((16Lk^3), 2\sigma^2, (\frac{ck}{4L})^3), \text{ and } \eta_t = k/(w + \sigma^2)$
e constants	. Assume	$\mathbb{E}[\ \mathscr{G}_{\boldsymbol{\eta}_{t^*}}(\mathbf{w}_{t_*})\ ^2] \leq \widetilde{O}\left($
		where t_* is sampled from $\{1, \ldots, T\}$.
schitz.	1 • .	
ny $\mathbf{w}_1, \mathbf{w}_2$ ai	nd its	Theorem RECOVER
W ₁ .		Assume that assumption 1, 2 hold and define co
t $\leq \sigma_{o'}^2.$		setting $\eta_k = \min\{\frac{\sqrt{\mu \varepsilon_k L}}{2c\sigma}, \frac{1}{16L}\}, T_k = O(\max\{\frac{1}{\mu^3}, K = O(\log(\varepsilon_1/\varepsilon))\})$ stages, the output of RECO
0		

[3] Yang, J., Kiyavash, N., He, N. (2020). Global convergence and variance-reduced optimization for a class of nonconvex-nonconcave minimax problems. [4] Guo, Z., Yuan, Z., Yan, Y., Yang, T. (2020). Fast Objective Duality Gap Convergence for Nonconvex-Strongly-Concave Min-Max Problems. [5] Wang, M., Liu, J., Fang, E. X. (2017). Accelerating stochastic composition optimization. Journal of Machine Learning Research. [6] Zhang, J., Xiao, L. (2019). A stochastic composite gradient method with incremental variance reduction.

Advances in Neural Information Processing Systems

(3)

ns: COVER & RECOVER

 $\eta \nabla g(\mathbf{w})^{\top} \nabla f(g(\mathbf{w}))).$

e standard gradient measure, i.e., $\mathscr{G}_n(\mathbf{w}) = \mathbf{w}$

is defined as the number of samples z in $\eta > 0$ or $\mathbb{E}[F(\mathbf{w}) - F_*] \leq \varepsilon$. False)

tes: $u_1 = g_z(\mathbf{w}_1), \ \mathbf{v}_1 = \nabla g_z(\mathbf{w}_1)$

$$-a_{t+1})(u_t - g_{\mathbf{z}_{t+1}}(\mathbf{w}_t))$$
$$-a_{t+1})(\mathbf{v}_t - \nabla g_{\mathbf{z}_{t+1}}(\mathbf{w}_t))$$

 $au \in \{1,\ldots,T\}.$

nstruct the estimates $u_0 = g_{\mathbf{z}_0}(\mathbf{w}_0), \ \mathbf{v}_0 =$

 T_k , True) 'ER

(Non-convex)

$$k = \frac{C\sigma^{2/3}}{L}, \ c = 128L + \sigma^2/(7Lk^3), \ w =$$

$${}^{2}t)^{1/3}.$$
 The output of COVER satisfies
$$\left(\frac{\Delta_F}{T^{2/3}} + \frac{\sigma^2}{T^{2/3}}\right).$$
(4)

(PL Condition)

constants $\varepsilon_1 = \frac{c^2 \sigma^2}{64 \mu L^4}$ and $\varepsilon_k = \varepsilon_1 / 2^{k-1}$. By $\frac{96c\sigma}{3/2\sqrt{\epsilon_k}L}, \frac{2c^2\sigma^2}{\mu L^2\epsilon_k}, \frac{\Delta_F}{\sigma^2}\}), c = 104L^2$, then after DVER satisfies $\mathbb{E}[F(\mathbf{w}_K) - F_*] \leq \varepsilon$.

Summary of properties of state-of-the-art algorithms for solving our DRO problem. The sample complexity is measured in terms of finding an ε -stationary point w/o PL condition, i.e., $\|\nabla F(\mathbf{w})\|^2 \leq \varepsilon$, or achieving ε -objective gap, i.e., $F(\mathbf{w}) - \min_{\mathbf{w}} F(\mathbf{w}) \leq \varepsilon$ with PL condition. O omits a logarithmic dependence over ε . n represents the size of datasets for a finite sum problem, d denotes the dimension of w. GDS represents whether the step size is geometrically decreased. Sample Complexity batch size GDS η Memory Algorithms Settings



Emperical Studies on Multi-class Imbalance Deep Learning Tasks

Datasets:



RECOVER optimized DRO vs SGD optimized ERM

We compare the test accuracy learned by optimizing DRO using RECOVER and optimizing ERM using SGD on the imbalanced datasets: STL10, CIFAR10, CIFAR100, with four imbalance ratio $\rho = \{0.02, 0.05, 0.1, 0.2\}.$ Table: Test accuracy (%), mean (std), of SGD for ERM and RECOVER for DRO over 5 independent runs.

IMRATIO	STL10		CIFA	AR10	CIFAR100	
	SGD	RECOVER	SGD	RECOVER	SGD	RECOVER
0.02	37.97 (0.78)	38.08 (0.59)	65.36(0.64)	66.14 (0.48)	38.99 (0.62)	39.45 (0.56)
0.05	41.12 (0.94)	42.68 (0.60)	74.74 (0.71)	75.90 (0.33)	45.79 (0.69)	44.47 (0.66)
0.1	46.03 (0.96)	48.94 (0.86)	79.32 (0.42)	80.93 (0.31)	49.45 (0.5)	50.84 (0.86)
0.2	51.75 (1.14)	56.06 (1.26)	84.84 (0.51)	85.93 (0.14)	55.80 (0.74)	56.90 (0.42)

Efficient as Fine-tune Methods

Dataset: ImageNet-LT, Places-LT



different methods.

Code can be found at: https://github.com/qiqi-helloworld/RECOVER

$\overline{O(n/\varepsilon+1/\varepsilon^2)}$	<i>O</i> (1)	X	O(n+d)	Primal-Dual
$O(1/\varepsilon^2)$	O(1)	Х	O(d)	Compositional
$O(1/\epsilon^{3/2})$	$O(1/\varepsilon)$	Х	O(d)	Compositional
$\widetilde{O}(1/arepsilon^{3/2})$	O(1)	Х	O(d)	Compositional
$O(1/\mu^2 \varepsilon)$	<i>O</i> (1)	Х	O(n+d)	Primal-Dual
$O(1/\mu^2 \varepsilon)$	<i>O</i> (1)	\checkmark	O(n+d)	Primal-Dual
$\widetilde{O}(1/\muarepsilon)$	$O(1/\varepsilon)$	Х	O(d)	Compositional
$O(1/\mu \epsilon)$	<i>O</i> (1)	\checkmark	O(d)	Compositional
		_		_

Half	Last Half	batch	Classes	Size	Network Arch
	500	32	10	5000	Resnet20

tical Imbalanced Datatset		64	1010	265,213	Inception-V3
	500	128	100	50000	Resnet20
	5000	128	10	50000	Resnet20
	500	32	10	5000	Resnet20

Models: ResNet50, ResNet152

	Model	ImageNet-LT	Places-LT
	Pretrained	40.50	23.28
	CE (SGD)	41.29 (3e-3)	27.47 (1e-3)
	Focal (SGD)	41.10 (2e-2)	27.64 (6e-3)
100	DRO (RECOVER)	42.30 (4e-4)	28.75 (4e-5)

Figure 1:Left: Test Accuracy vs λ on CIFAR10 data; Right: Test accuracy (%) of finetuned models by

^[1] Rafique, H., Liu, M., Lin, Q., Yang, T. (1810). Non-convex min-max optimization: provable algorithms and applications in machine learning (2018). [2] Yan, Y., Xu, Y., Lin, Q., Liu, W., Yang, T. (2020). Sharp analysis of epoch stochastic gradient descent

ascent methods for min-max optimization