

Finite-Sum Coupled Compositional Stochastic Optimization: Theory and Applications

Bokun Wang and Tianbao Yang
bokun-wang.github.io github.com/bokun-wang/SOX



IOWA

The Problem

Minimize a sum of compositions: inner-level function of each summand **coupled** with outer index.

$$\min_{\mathbf{w} \in \Omega} F(\mathbf{w}),$$

$$F(\mathbf{w}) := \frac{1}{n} \sum_{\mathbf{z}_i \in \mathcal{D}} f_i(g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}_i))$$

Coupled

Finite-Sum Coupled Compositional Optimization (FCCO)

Applications

1. Average precision maximization

$$F(\mathbf{w}) = -\frac{1}{|\mathcal{S}_+|} \sum_{\mathbf{z}_i \in \mathcal{S}_+} \sum_{\mathbf{z} \in \mathcal{S}} \ell(h_{\mathbf{w}}(\mathbf{z}) - h_{\mathbf{w}}(\mathbf{z}_i)).$$

$h_{\mathbf{w}}(\mathbf{z})$: model prediction; \mathcal{S}_+ : positive data; \mathcal{S} : whole data; $\ell(\cdot)$: monotonically decreasing upper bound on the 0-1 loss (e.g. hinge loss).

$$f_i(\cdot) = -\frac{[\cdot]_1}{[\cdot]_2}, \quad g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}) = \left[\sum_{\mathbf{z} \in \mathcal{S}_+} \ell(h_{\mathbf{w}}(\mathbf{z}) - h_{\mathbf{w}}(\mathbf{z}_i)) \right].$$

2. Ranking by p -norm push

$$F(\mathbf{w}) = \sum_{\mathbf{z}_i \in \mathcal{S}_-} \left(\sum_{\mathbf{z}_j \in \mathcal{S}_+} \ell(h_{\mathbf{w}}(\mathbf{z}_j) - h_{\mathbf{w}}(\mathbf{z}_i)) \right)^p.$$

$p > 1$: order; \mathcal{S}_- : negative data; others same as AP maximization.

$$f_i(\cdot) = (\cdot)^p, \quad g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}_+) = \sum_{\mathbf{z}_j \in \mathcal{S}_+} \ell(h_{\mathbf{w}}(\mathbf{z}_j) - h_{\mathbf{w}}(\mathbf{z}_i)).$$

3. Neighborhood component analysis

$$F(A) = -\sum_{\mathbf{z}_i \in \mathcal{D}} \frac{\sum_{\mathbf{z} \in \mathcal{C}_i} \exp(-\|A\mathbf{z}_i - A\mathbf{z}\|^2)}{\sum_{\mathbf{z} \in \mathcal{S}_i} \exp(-\|A\mathbf{z}_i - A\mathbf{z}\|^2)}.$$

A : embedding matrix; \mathcal{D} : data set; \mathcal{C}_i : nearest-neighbour set of \mathbf{z}_i ; \mathcal{S}_i : leave-one-out set of \mathbf{z}_i .

$$f_i(\cdot) = -\frac{[\cdot]_1}{[\cdot]_2}, \quad g(A; \mathbf{z}_i, \mathcal{S}_i) = \left[\sum_{\mathbf{z} \in \mathcal{C}_i} \exp(-\|A\mathbf{z}_i - A\mathbf{z}\|^2) \right].$$

4. Listwise Ranking, e.g. ListNet

$$F(\mathbf{w}) = -\sum_q \sum_{\mathbf{x}_i^q \in \mathcal{S}_q} P(y_i^q) \log \frac{\exp(h_{\mathbf{w}}(\mathbf{x}_i^q; \mathbf{q}))}{\sum_{\mathbf{x} \in \mathcal{S}_q} \exp(h_{\mathbf{w}}(\mathbf{x}; \mathbf{q}))}.$$

$h_{\mathbf{w}}(\mathbf{x}_i^q; \mathbf{q})$: prediction score of item \mathbf{x}_i^q ; $\mathcal{Q} = \{\mathbf{q}_1, \dots, \mathbf{q}_n\}$: queries; $\mathcal{S}_q = \{(\mathbf{x}_1^q, y_1^q), \dots, (\mathbf{x}_n^q, y_n^q)\}$: item-relevance pairs.

$$f_i(\cdot) = \log(\cdot), \quad \mathcal{D} = \{(\mathbf{q}, \mathbf{x}_i^q) \mid P(y_i^q) > 0\},$$

$$g(\mathbf{w}; \mathbf{x}_i^q, \mathcal{S}_q) = \sum_{\mathbf{x} \in \mathcal{S}_q} \exp(h_{\mathbf{w}}(\mathbf{x}; \mathbf{q}) - h_{\mathbf{w}}(\mathbf{x}_i^q; \mathbf{q})).$$

Applications (cont'd)

5. Survival Analysis

$$F(\mathbf{w}) = \sum_{i: E_i=1} \log \left(\sum_{j \in \mathcal{S}(T_i)} \exp(h_{\mathbf{w}}(\mathbf{z}_j) - h_{\mathbf{w}}(\mathbf{z}_i)) \right).$$

\mathbf{z}_i : patient feature; $h_{\mathbf{w}}(\mathbf{z}_i)$: predicted risk; $E_i = 1$: observable event of interest (e.g., death); T_i : time interval between data collection and event; $\mathcal{S}(t) = \{i : T_i \geq t\}$: patients at risk of failure at time t .

$$f_i(\cdot) = -\log(\cdot), \quad \mathcal{D} = \{\mathbf{z}_i \mid E_i = 1\},$$

$$g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}(T_i)) = \sum_{j \in \mathcal{S}(T_i)} \exp(h_{\mathbf{w}}(\mathbf{z}_j) - h_{\mathbf{w}}(\mathbf{z}_i)).$$

6. Latent Variable Model

$$F(\mathbf{w}) = -\sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \log \sum_{\mathbf{s} \in \mathcal{S}} \Pr(y_i \mid \mathbf{s}, \mathbf{x}_i) \Pr(\mathbf{s} \mid \mathbf{x}_i).$$

\mathbf{s} : discrete latent variable; \mathcal{S} : support set of the latent variable \mathbf{s} .

$$f_i(\cdot) = -\log(\cdot), \quad \mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}.$$

$$g(\mathbf{w}; (\mathbf{s}_i, y_i), \mathcal{S}) = \sum_{\mathbf{s} \in \mathcal{S}} \Pr(y_i \mid \mathbf{s}, \mathbf{x}_i) \Pr(\mathbf{s} \mid \mathbf{x}_i).$$

Why FCCO?

• Finite-sum optimization (FO) problem:

$$F(\mathbf{w}) = \frac{1}{n} \sum_{\mathbf{z}_i \in \mathcal{D}} \ell(\mathbf{w}; \mathbf{z}_i).$$

Algorithms for FO: SGD, SGDM, Adam...

Why not view applications 1.-6. as FO?

⇒ Computing an **unbiased** stochastic estimator of $\nabla F(\mathbf{w})$ is expensive when inner batch \mathcal{S}_i is large!

• Finite-sum compositional optimization (FCO):

$$F(\mathbf{w}) = \frac{1}{n} \sum_{\mathbf{z}_i \in \mathcal{D}} f_i(g(\mathbf{w}; \mathcal{S})).$$

Algorithms for FCO: SCGD, NASA, ALSET...

We can reformulate an FCCO problem as FCO problem by re-writing $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \hat{f}_i(\mathbf{g}(\mathbf{w}; \mathcal{S}))$, where $\hat{f}_i(\cdot) = f_i(\mathbb{I}_i)$, $\mathbb{I}_i = [0_{d \times d}, \dots, I_{d \times d}, \dots, 0_{d \times d}]$, $\mathbf{g}(\mathbf{w}; \mathcal{S}) = [g(\mathbf{w}; \mathbf{z}_1, \mathcal{S}_1)^\top, \dots, g(\mathbf{w}; \mathbf{z}_n, \mathcal{S}_n)^\top]^\top$, $\mathcal{S} = \mathcal{S}_1 \cup \dots \cup \mathcal{S}_i \cup \dots \cup \mathcal{S}_n$.

Why not view applications 1.-6. as FCO?

⇒ FCO algorithms compute stochastic estimators for **all** n components of $\mathbf{g}(\mathbf{w}; \mathcal{S})$ (inefficient)!

Why FCCO? (cont'd)

• Conditional stochastic optimization (CSO):

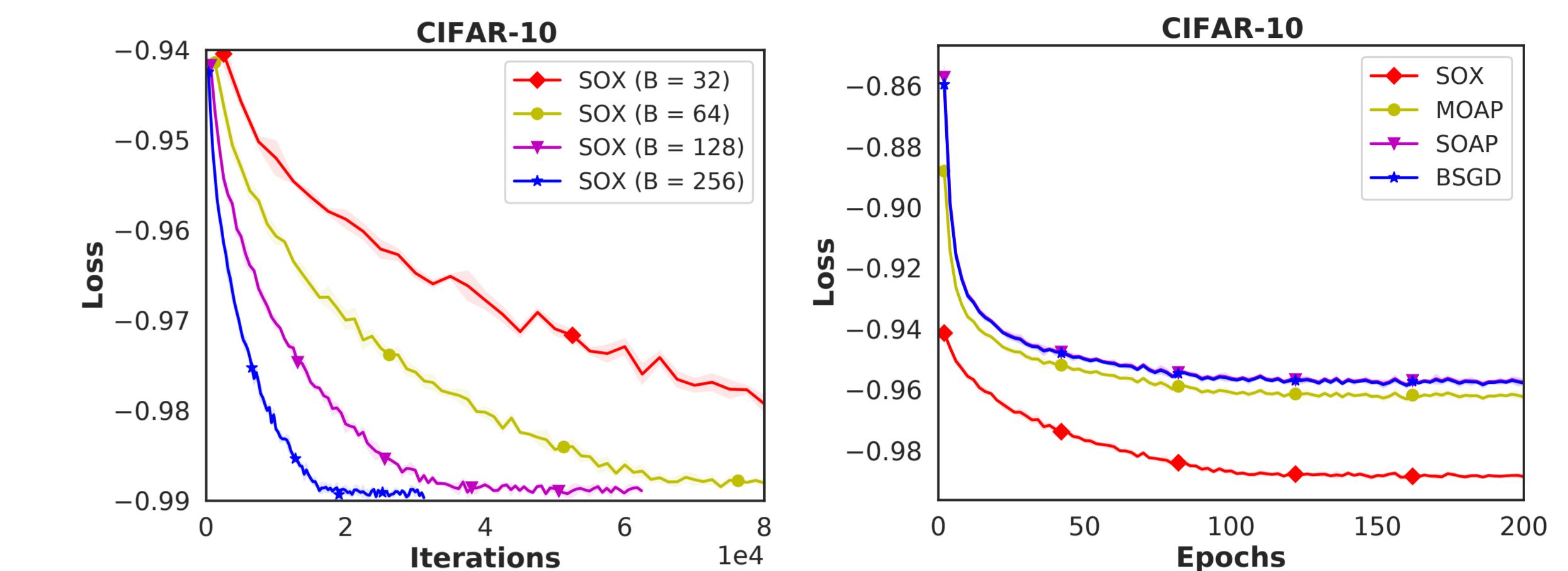
$$\mathbb{E}_{\xi} f_{\xi} \left(\mathbb{E}_{\zeta | \xi} [g_{\zeta}(\mathbf{w}; \xi)] \right).$$

Algorithms for CSO: BSGD, $O(\epsilon^{-2})$ batch size. FCCO is a special case of CSO whose outer problem has finite support.

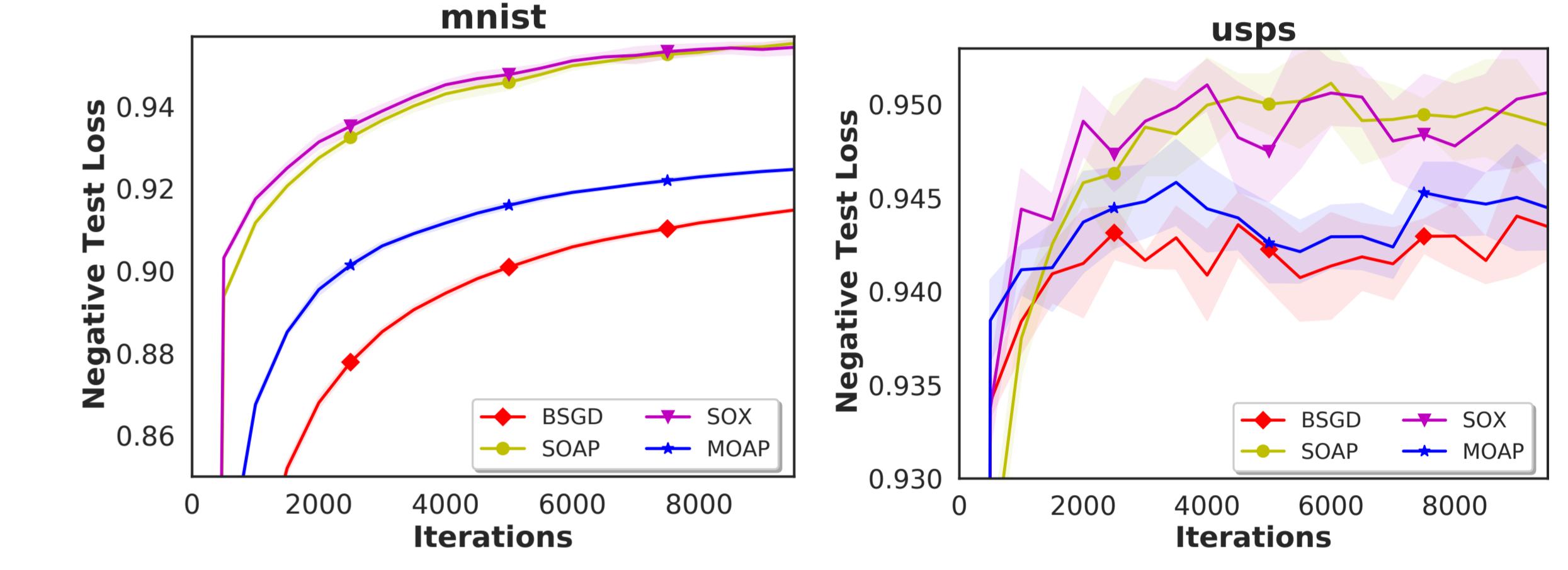
Why not view applications 1.-6. as CSO?

⇒ CSO algorithm has worse sample complexity because the finite-sum structure is not used!

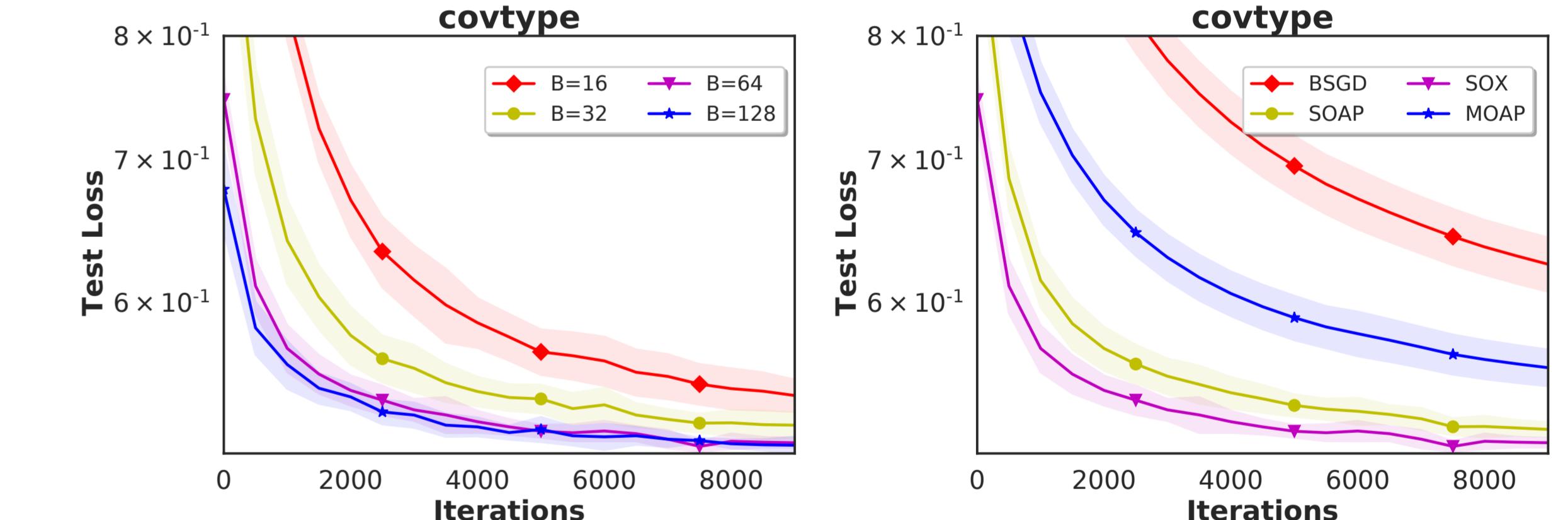
Numerical Experiments (cont'd)



• Neighborhood component analysis



• Ranking by p -norm push



BS-PnP: boosting-style baseline.

	BS-PnP	BSGD	SOX
Test Loss (\downarrow)	0.778	0.625 \pm 0.018	0.516 \pm 0.003
Time (s) (\downarrow)	6043.90	4.20 \pm 0.08	4.62 \pm 0.10
		ijcnn1	
Algorithms	BS-PnP	BSGD	SOX
Test Loss (\downarrow)	0.268	0.202 \pm 0.001	0.128 \pm 0.002
Time (s) (\downarrow)	648.06	4.02 \pm 0.04	4.15 \pm 0.06

Numerical Experiments

Iteration Complexity

• Strongly convex: $O\left(\frac{n\epsilon^{-1}}{\mu^2 B_1 B_2}\right)$.

• Convex: $O\left(\frac{n\epsilon^{-3}}{B_1 B_2}\right)$.

• Convex and monotone: $O\left(\frac{n\epsilon^{-2}}{B_1}\right)$.

• Non-convex: $O\left(\frac{n\epsilon^{-4}}{B_1 B_2}\right)$.

References

- [SCGD] Wang, M., Fang, E. X., and Liu, H. *Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions*. *Mathematical Programming*, 161(1), 419-449.
- [NASA] Ghadimi, S., Ruszczynski, A., and Wang, M. *A single timescale stochastic approximation method for nested stochastic optimization*. *SIAM Journal on Optimization*, 30(1), 960-979.
- [BSGD] Hu, Y., Zhang, S., Chen, X., and He, N. *Biased Stochastic First-Order Methods for Conditional Stochastic Optimization and Applications in Meta Learning*. *Advances in Neural Information Processing Systems*, 33, 2759-2770.
- [SOAP] Qi, Q., Luo, Y., Xu, Z., Ji, S., and Yang, T. *Stochastic optimization of areas under precision-recall curves with provable convergence*. *Advances in Neural Information Processing Systems*, 34, 1752-1765.
- [MOAP] Wang, G., Yang, M., Zhang, L., and Yang, T. *Momentum accelerates the convergence of stochastic aupre maximization*. *International Conference on Artificial Intelligence and Statistics*, 25, 3753-3771.