Finite-Sum Coupled Compositional Stochastic Optimization Theory and Applications

Bokun Wang and Tianbao Yang





Empirical Risk Minimization (ERM)

$$\hat{R}(h) = rac{1}{n} \sum_{i=1}^{n} L(h(\mathbf{x}_i), y_i).$$
 $\mathcal{D} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$
Sample Loss Feature Label $\hat{h} = rgmin\hat{R}(h)$
 $h \in \mathcal{H}$

`

$$egin{aligned} &\min_{h\in\mathcal{H}}\hat{R}(h), \ \hat{R}(h) = rac{1}{n}\sum_{i=1}^n L(h(\mathbf{x}_i),y_i). \end{aligned}$$
 Hypothesis parameterized by \mathbf{w} $&\min_{\mathbf{w}\in\Omega}F(\mathbf{w}), \quad F(\mathbf{w}) := rac{1}{n}\sum_{\mathbf{z}_i\in\mathcal{D}}\ell(\mathbf{w};\mathbf{z}_i) \end{aligned}$

Gradient Descent

 $n = |\mathcal{D}|$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta rac{1}{n} \sum_{\mathbf{z}_i \in \mathcal{D}}
abla \ell(\mathbf{w}; \mathbf{z}_i)$$

$$\begin{split} \min_{h \in \mathcal{H}} \hat{R}(h), \ \hat{R}(h) &= \frac{1}{n} \sum_{i=1}^{n} L(h(\mathbf{x}_{i}), y_{i}). \\ \text{Hypothesis parameterized by } \mathbf{w} \\ \min_{\mathbf{w} \in \Omega} F(\mathbf{w}), \ F(\mathbf{w}) &:= \frac{1}{n} \sum_{\mathbf{z}_{i} \in \mathcal{D}} \ell(\mathbf{w}; \mathbf{z}_{i}) \\ \hline \frac{1}{n} \sum_{\mathbf{z}_{i} \in \mathcal{D}} \nabla \ell(\mathbf{w}; \mathbf{z}_{i}) \\ Expensive when n is \\ large ! \end{split}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta
abla F(\mathbf{w})$$

Gradient Descent

$$egin{aligned} &\min_{h\in\mathcal{H}}\hat{R}(h), \;\; \hat{R}(h) = rac{1}{n}\sum_{i=1}^n L(h(\mathbf{x}_i),y_i). \end{aligned}$$
 Hypothesis barameterized by \mathbf{w} $&\min_{\mathbf{w}\in\Omega}F(\mathbf{w}), \;\;\; F(\mathbf{w}) := rac{1}{n}\sum_{\mathbf{z}_i\in\mathcal{D}}\ell(\mathbf{w};\mathbf{z}_i) \end{aligned}$

Stochastic Gradient Descent

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \hat{
abla} F(\mathbf{w})$$

 $\mathbb{E}[\hat{
abla}F(\mathbf{w})] =
abla F(\mathbf{w})$

$$egin{aligned} &\min_{h\in\mathcal{H}}\hat{R}(h), \;\; \hat{R}(h) = rac{1}{n}\sum_{i=1}^n L(h(\mathbf{x}_i),y_i). \end{aligned}$$
 Hypothesis barameterized by \mathbf{w} $&\min_{\mathbf{w}\in\Omega}F(\mathbf{w}), \;\;\; F(\mathbf{w}) := rac{1}{n}\sum_{\mathbf{z}_i\in\mathcal{D}}\ell(\mathbf{w};\mathbf{z}_i) \end{aligned}$

Stochastic Gradient Descent

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \hat{
abla} F(\mathbf{w})$$

Unbiased estimator, e.g., $abla \ell(\mathbf{w}; \mathbf{z}_i)$

Independent of n. Looks good?

Surrogate of Average Precision (AP) Maximization

$$\begin{split} F(\mathbf{w}) &= -\frac{1}{|\mathcal{S}_{+}|} \sum_{\mathbf{x}_{i} \in \mathcal{S}_{+}} \frac{\sum_{\mathbf{x} \in \mathcal{S}_{+}} \ell(h_{\mathbf{w}}(\mathbf{x}) - h_{\mathbf{w}}(\mathbf{x}_{i}))}{\sum_{\mathbf{x} \in \mathcal{S}} \ell(h_{\mathbf{w}}(\mathbf{x}) - h_{\mathbf{w}}(\mathbf{x}_{i}))} \\ & \\ \text{Positive} \qquad & \\ \text{Data} \qquad & \mathcal{S} = \mathcal{S}_{+} \cup \mathcal{S}_{-} \end{split}$$

Surrogate of $\min_{\mathbf{w}\in\Omega}F(\mathbf{w}), \quad F(\mathbf{w}):=rac{1}{n}\sum_{\mathbf{z}_i\in\mathcal{D}}\ell(\mathbf{w};\mathbf{z}_i)$ Average Precision (AP) Maximization

$$F(\mathbf{w}) = -rac{1}{|\mathcal{S}_+|} \sum_{\mathbf{x}_i \in \mathcal{S}_+} egin{pmatrix} \sum {\sum_{\mathbf{x} \in \mathcal{S}_+} \ell(h_{\mathbf{w}}(\mathbf{x}) - h_{\mathbf{w}}(\mathbf{x}_i)) \over \sum_{\mathbf{x} \in \mathcal{S}} \ell(h_{\mathbf{w}}(\mathbf{x}) - h_{\mathbf{w}}(\mathbf{x}_i))} \end{pmatrix}} \ell(\mathbf{w}; \mathbf{z}_i)$$

Stochastic Gradient Descent

$$\mathbf{w} \leftarrow \mathbf{w} - \eta
abla \ell(\mathbf{w}; \mathbf{z}_i)$$

Unbiased estimator is still expensive !



$$\mathbf{w} \leftarrow \mathbf{w} - \eta
abla \ell(\mathbf{w}; \mathbf{z}_i)$$
 Infeasible !

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

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

How is it related to finite-sum optimization?

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

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

Take into account the cost of S_i

Finite-Sum Optimization (FO)

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

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

$F(\mathbf{w}) := \frac{1}{n} \sum_{\mathbf{z}_i \in \mathcal{D}} f_i(g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}_i))$ Finite-Sum Coupled Composition Optimization (FCCO)

- Bipartite ranking by p-norm Push $F(\mathbf{w}) = \frac{1}{|\mathcal{S}_{-}|} \sum_{\mathbf{z}_{i} \in \mathcal{S}_{-}} \left(\frac{1}{|\mathcal{S}_{+}|} \sum_{\mathbf{z}_{j} \in \mathcal{S}_{+}} \ell(h_{\mathbf{w}}(\mathbf{z}_{j}) - h_{\mathbf{w}}(\mathbf{z}_{i})) \right)^{p}$
 - Neighborhood Component Analysis

$$F(A) = -\sum_{\mathbf{x}_i \in \mathcal{D}} \frac{\sum_{\mathbf{x} \in \mathcal{C}_i} \exp(-\|A\mathbf{x}_i - A\mathbf{x}\|^2)}{\sum_{\mathbf{x} \in \mathcal{S}_i} \exp(-\|A\mathbf{x}_i - A\mathbf{x}\|^2)} \mathcal{S}_i = \mathcal{D} \smallsetminus \{\mathbf{x}_i\} \qquad F(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n \left\|\mathbf{x}_i^\top \mathbf{w} - y_i\right\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 = \mathcal{C}_i = \{\mathbf{x}_j \in \mathcal{D} : y_j = y_i\}$$

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

Finite-Sum
Optimization (FO)

• Logistic regression

$$F(\mathbf{w}) = rac{1}{n} \sum_{i=1}^n \ln \Bigl(1 + e^{-y_i \langle \mathbf{w}, \mathbf{x}_i
angle} \Bigr) \, .$$

• Ridge regression

Stochastic Alg. for FCCO problems

Stochastic Alg. for FO problems

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

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

Stochastic Gradient (**Biased**); Sample both D and S_i $\min_{\mathbf{w}\in\Omega}F(\mathbf{w}),$

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

Stochastic Gradient (Unbiased); Sample ${\cal D}$

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

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

Wait ! We have already seen something similar ...

Hu et al. "Biased stochastic first-order methods for conditional stochastic optimization and applications in meta learning." NeurIPS 2020.

Conditional Stochastic Optimization (CSO)

 $\begin{array}{lll} \text{Goal: better sample} & \textit{Special Case:} & \text{BSGD, BSpiderBoost:} \\ \text{complexity \& O(1)} & \textit{Outer problem has} & O(\sqrt{T}) \text{ batch size} \\ \text{batch size !} & \textit{finite support} \\ F(\mathbf{w}) := \frac{1}{n} \sum_{\mathbf{z}_i \in \mathcal{D}} f_i(g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}_i)) & F(\mathbf{w}) = \mathbb{E}_{\xi} f_{\xi} \big(\mathbb{E}_{\zeta \mid \xi}[g_{\zeta}(\mathbf{w}; \xi)] \big) \\ \end{array}$

Wang et al. "Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions." Math. Program. 161(1-2):419–449, 2017.

Finite-Sum Composition Optimization (FCO)

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

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

 $\min_{\mathbf{w}\in\Omega}F(\mathbf{w}),$ $F(\mathbf{w}) = rac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}} f_i(g(\mathbf{w}; \mathcal{S}))$

Wang et al. "Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions." Math. Program. 161(1-2):419–449, 2017.

Finite-Sum Composition Optimization (FCO)

The NASA Algorithm for FCO problem

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

Sample mini-batches $\mathcal{B}_1 \subset \mathcal{D}, \mathcal{B}_2 \subset \mathcal{S}$

$$egin{aligned} & u \leftarrow (1-\gamma)u + \gamma g(\mathbf{w};\mathcal{B}_2) \ & \mathbf{v} \leftarrow (1-eta)\mathbf{v} + eta rac{1}{|\mathcal{B}_1|} \sum_{\mathbf{z}_i \in \mathcal{B}_1}
abla g(\mathbf{w};\mathcal{B}_2)
abla f_i(u) \end{aligned}$$

 $\mathbf{w} \longleftarrow \mathbf{w} - \eta \mathbf{v}$

Apply NASA to FCCO? $F(\mathbf{w}) := rac{1}{n} \sum_{i \in \mathcal{T}} f_i(g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}_i))$ $\mathbf{g}(\mathbf{w};\mathcal{B}_2) = \left[g(\mathbf{w};\mathbf{z}_1,\mathcal{B}_{2,1})^ op,\ldots,g(\mathbf{w};\mathbf{z}_n,\mathcal{B}_{2,n})^ op
ight]^+$ Reformulation $F(\mathbf{w}) = rac{1}{n} \sum_{i=1}^n \hat{f}_i(\mathbf{g}(\mathbf{w};\mathcal{S}))$ Each iteration: sample and update $\mathbf{g}(\mathbf{w};\mathcal{S}) = \left[g(\mathbf{w};\mathbf{z}_1,\mathcal{S}_1)^ op,\ldots,g(\mathbf{w};\mathbf{z}_n,\mathcal{S}_n)^ op
ight]^ op$ for all n coordinates ! $\mathcal{S} = \mathcal{S}_1 \cup \cdots \mathcal{S}_i \cdots \cup \mathcal{S}_n$ Not efficient when n is large. $\hat{f}_i(\cdot) = f_i(\mathbb{I}_i \cdot) \quad \mathbb{I}_i := [0_{d imes d}, \dots, I_{d imes d}, \dots, 0_{d imes d}]$

Say n = 5,
$$B_1 = 2$$

Remedy: NASA + Rand. Sparsification

 $\mathbf{g}(\mathbf{w}; \mathcal{B}_2) = \left[g(\mathbf{w}; \mathbf{z}_1, \mathcal{B}_{2,1})^\top, g(\mathbf{w}; \mathbf{z}_2, \mathcal{B}_{2,2})^\top, g(\mathbf{w}; \mathbf{z}_3, \mathcal{B}_{2,3})^\top, g(\mathbf{w}; \mathbf{z}_4, \mathcal{B}_{2,4})^\top, g(\mathbf{w}; \mathbf{z}_5, \mathcal{B}_{2,5})^\top\right]$

$$\mathbf{g}(\mathbf{w};\mathcal{B}_2) = \left[0,g(\mathbf{w};\mathbf{z}_2,\mathcal{B}_{2,2})^ op,0,0,g(\mathbf{w};\mathbf{z}_5,\mathcal{B}_{2,5})^ op
ight]^ op \left[rac{n}{B_1}
ight]$$

Only compute $B_1 \ll n$ coordinates.

Randomly replace others with zeros

$$u (1-\gamma)u + \gamma g(\mathbf{w}; \mathcal{B}_2)) \hspace{0.2cm} u \in \mathbb{R}^n$$

1) overflow? 2) per-iteration cost of rescaling (n-B₁) coordinates by (1 - γ). 3) no speed-up w.r.t. B₂ (Proposition

3.5 in Khirirat et al. 2018).

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

Sample mini-batches
$$\mathcal{B}_{1}^{t} \subset \mathcal{D}, \mathcal{B}_{i,2}^{t} \subset \mathcal{S}_{i}$$

 $u_{i}^{t} = \begin{cases} (1 - \gamma)u_{i}^{t-1} + \gamma g(\mathbf{w}^{t}; \mathbf{z}_{i}, \mathcal{B}_{i,2}^{t}), & \mathbf{z}_{i} \in \mathcal{B}_{1}^{t} \\ u_{i}^{t-1}, & \mathbf{z}_{i} \notin \mathcal{B}_{1}^{t} \end{cases}$
Only update and sample for a subset of coordinates !
 $\mathbf{v}^{t} = (1 - \beta)\mathbf{v}^{t-1} + \beta \frac{1}{B_{1}} \sum_{\mathbf{z}_{i} \in \mathcal{B}_{1}^{t}} \nabla g(\mathbf{w}^{t}; \mathbf{z}_{i}, \mathcal{B}_{i,2}^{t}) \nabla f_{i}(u_{i}^{t-1})$
 $\mathbf{w}^{t+1} = \mathbf{w}^{t} - \eta_{t} \mathbf{v}^{t}$
Per-iteration computation cost: O(B₁)

Finite-Sum Coupled Composition **Optimization** (FCCO) (NEW) The SOX Algorithm $F(\mathbf{w}) := rac{1}{n} \sum_{\mathbf{z}_i \in \mathcal{D}} f_i(g(\mathbf{w}; \mathbf{z}_i, \mathcal{S}_i))$ $u_i^t = egin{cases} (1-\gamma)u_i^{t-1} + \gamma gig(\mathbf{w}^t; \mathbf{z}_i, \mathcal{B}_{i,2}^tig), & \mathbf{z}_i \in \mathcal{B}_1^t \ u_i^{t-1}, & \mathbf{z}_i
ot\in \mathcal{B}_1^t \end{cases}$ $u_i^t = egin{cases} u_i^{t-1} - \gammaig(u_i^{t-1} - gig(\mathbf{w}^t; \mathbf{z}_i, \mathcal{B}_{i,2}^tig)ig), & \mathbf{z}_i \in \mathcal{B}_1^t \ u_i^t, & \mathbf{z}_i
otin \mathcal{B}_1^t \end{cases}$ Stochastic block coordinate descent $\min_{\mathbf{u} = \left[u_1, \ldots, u_n
ight]^ op} rac{\mathbf{1}}{2} \sum_{\mathbf{r}_i \in \mathcal{D}} \left\|u_i - gig(\mathbf{w}^t; \mathbf{z}_i, \mathcal{S}_iig)
ight\|^2$

Finite-Sum <u>Coupled</u> Composition Optimization (<u>FCCO</u>)

(NEW) The SOX Algorithm

Sample mini-batches $\mathcal{B}_1^t \subset \mathcal{D}, \mathcal{B}_{i,2}^t \subset \mathcal{S}_i$

$$u_i^t = egin{cases} (1-\gamma)u_i^{t-1} + \gamma gig(\mathbf{w}^t; \mathbf{z}_i, \mathcal{B}_{i,2}^tig), & \mathbf{z}_i \in \mathcal{B}_1^t \ u_i^{t-1}, & \mathbf{z}_i
ot\in \mathcal{B}_1^t \end{cases}$$

$$\mathbf{v}^t = (1-eta)\mathbf{v}^{t-1} + eta rac{1}{B_1} \sum_{\mathbf{z}_i \in \mathcal{B}_1^t}
abla g(\mathbf{w}^t; \mathbf{z}_i, \mathcal{B}_{i,2}^t)
abla f_i(u_i^{t-1}) \quad u_i^t \text{ is more intuitive ?}$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \mathbf{v}^t$$



Originally proposed for AP maximization * extra assumption: monotonicity

Bipartite Ranking
by p-norm Push $F(\mathbf{w}) = \frac{1}{|\mathcal{S}_{-}|} \sum_{\mathbf{z}_{i} \in \mathcal{S}_{-}} \left(\frac{1}{|\mathcal{S}_{+}|} \sum_{\mathbf{z}_{j} \in \mathcal{S}_{+}} \ell(h_{\mathbf{w}}(\mathbf{z}_{j}) - h_{\mathbf{w}}(\mathbf{z}_{i})) \right)^{p}$

A boosting-style deterministic algorithm

Algorithms	BS-PnP	SOX
Test Loss (\downarrow)	0.778	$\textbf{0.516} \pm \textbf{0.003}$
Time (s) (\downarrow)	6043.90	4.62 ± 0.10
Algorithms	BS-PnP	SOX
Test Loss (\downarrow)	0.268	$\textbf{0.128} \pm \textbf{0.002}$
Time (s) (\downarrow)	648.06	4.15 ± 0.06



Neighborhood Component Analysis

$$F(A) = -\sum_{\mathbf{x}_i \in \mathcal{D}} \frac{\sum_{\mathbf{x} \in \mathcal{C}_i} \exp(-\|A\mathbf{x}_i - A\mathbf{x}\|^2)}{\sum_{\mathbf{x} \in \mathcal{S}_i} \exp(-\|A\mathbf{x}_i - A\mathbf{x}\|^2)}$$
$$\mathcal{C}_i = \{\mathbf{x}_j \in \mathcal{D} : y_j = y_i\}$$
$$\mathcal{S}_i = \mathcal{D} \smallsetminus \{\mathbf{x}_i\}$$



More applications of SOX: partial AUC [Zhu et. al. 2022], NDCG [Qiu et.al. 2022], contrastive learning [Yuan et.al. 2022], listwise ranking, survival analysis, etc.

$\textbf{AP Maximization} \qquad F(\mathbf{w}) = -\frac{1}{|\mathcal{S}_+|} \sum_{\mathbf{x}_i \in \mathcal{S}_+} \frac{\sum_{\mathbf{x} \in \mathcal{S}_+} \ell(h_{\mathbf{w}}(\mathbf{x}) - h_{\mathbf{w}}(\mathbf{x}_i))}{\sum_{\mathbf{x} \in \mathcal{S}} \ell(h_{\mathbf{w}}(\mathbf{x}) - h_{\mathbf{w}}(\mathbf{x}_i))}$





Thank you !