# A Single Cut Proximal Bundle Method for Stochastic Convex Composite Optimization

Jiaming Liang

Yale University

October 16, 2022

Joint work with Vincent Guigues (FGV) and Renato Monteiro (Georgia Tech)

INFORMS Annual Meeting 2022, Indianapolis, IN

## Introduction

**Main problem**

$$\phi_* := \min \left\{ \phi(x) := f(x) + h(x) : x \in \mathbb{R}^n \right\}, \quad f(x) = \mathbb{E}_\xi[F(x, \xi)]$$

E.g., two-stage convex stochastic program

$$\min\{f_1(x) + \mathbb{E}[Q(x, \xi)] : x \in X\}$$

where $Q(x, \xi) = \min\{f_2(x, y, \xi) : g_2(x, y, \xi) \leq 0, y \in Y\}$.

An instance of the main problem with

$$h(x) = \delta_X(x), \quad F(x, \xi) = f_1(x) + Q(x, \xi).$$

**Goal**: SA-type algorithm based on the proximal bundle (PB) method

# Assumptions

**Stochastic convex composite optimization**

$$\phi_* := \min \left\{ \phi(x) := f(x) + h(x) : x \in \mathbb{R}^n \right\}, \quad f(x) = \mathbb{E}_\xi[F(x,\xi)]$$

**Black-box model**

(A1) $f$ is closed convex and $\operatorname{dom} f \supset \operatorname{dom} h$;

(A2) for almost every $\xi \in \Xi$, there exist a functional oracle $F(\cdot, \xi) : \operatorname{dom} h \to \mathbb{R}$ and a stochastic subgradient oracle $s(\cdot, \xi) : \operatorname{dom} h \to \mathbb{R}^n$ satisfying

$$f(x) = \mathbb{E}[F(x,\xi)], \quad f'(x) := \mathbb{E}[s(x,\xi)] \in \partial f(x);$$

(A3) for every $x \in \operatorname{dom} h$, we have $\mathbb{E}[\|s(x,\xi)\|^2] \le M^2$;

(A4) the set of optimal solutions $X^*$ is nonempty.

# Review of Deterministic PB

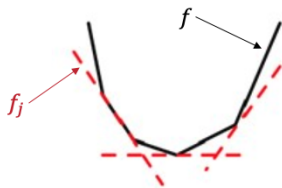Proximal point method: constructs a sequence of proximal problems.
E.g., Chambolle-Pock for saddle point, ADMM for distributed optimization.

Approximately solve the proximal problem by an iterative process

$$x^+ \leftarrow \min_{u \in \mathbb{R}^n} \left\{ f(u) + \frac{1}{2\lambda} \|u - x^c\|^2 \right\}.$$

Recursively build up a cutting-plane model

$$f_j(u) = \max_{0 \le i \le j-1} \{\ell_f(u; x_i) := f(x_i) + \langle f'(x_i), u - x_i \rangle\}$$

# Review of Deterministic PB

**Algorithm 1** PB (one cycle)

1. Construct a proximal problem

$$\min_{u \in \mathbb{R}^n} \left\{ f(u) + h(u) + \frac{1}{2\lambda} \|u - x^c\|^2 \right\};$$

2. **If** find an $(\varepsilon/2)$-solution to the current proximal problem, **then** change the prox-center; $\leftarrow$ serious

**Otherwise**, keep the prox-center, update the cutting-plane model and solve the prox subproblem based on the current model, i.e., $\leftarrow$ null

$$x_j = \operatorname*{argmin}_{u \in \mathbb{R}^n} \left\{ f_j(u) + \frac{1}{2\lambda} \|u - x^c\|^2 \right\}.$$

# Cutting-plane Model in the Stochastic Setting

A straightforward fact:

$$\mathbb{E}[\max\{X, Y\}] \geq \max\{\mathbb{E}[X], \mathbb{E}[Y]\}.$$

For a fixed $u$,

$$\mathbb{E}[\Gamma_j(u)] \geq \max_{0 \leq i \leq j-1}\{\mathbb{E}[F(x_i, \xi_i) + \langle s(x_i, \xi_i), u - x_i\rangle]\}.$$

On the other hand,

$$\max_{0 \leq i \leq j-1}\{\mathbb{E}[F(x_i, \xi_i) + \langle s(x_i, \xi_i), u - x_i\rangle]\}$$
$$= \max_{0 \leq i \leq j-1}\{f(x_i) + \langle f'(x_i), u - x_i\rangle\} \leq f(u)$$

So

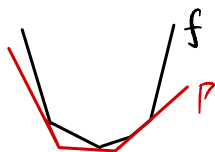$$\mathbb{E}[\Gamma_j(u)] \; ? \; f(u)$$

# Other bundle models

(E1) **single cut update**[1]: $\Gamma^+ = \Gamma_\tau^+ := \tau\Gamma + (1-\tau)\ell_f(\cdot;x)$.

(E2) **two cuts update:** assume $\Gamma = \max\{A_f, \ell_f(\cdot;x^-)\}$ where $A_f$ is an affine function satisfying $A_f \le f$, set

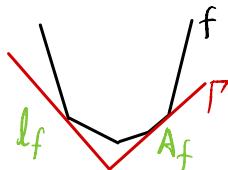$$\Gamma^+ = \max\{A_f^+, \ell_f(\cdot;x)\}$$
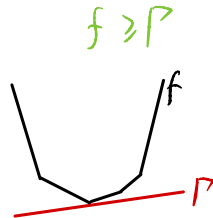
where $A_f^+ = \theta A_f + (1-\theta)\ell_f(\cdot;x^-)$.

Bundle of past information $\{(x_i, f(x_i), f'(x_i))\}$



Multiple cuts     Two cuts     One cut

---

[1]Liang and Monteiro, 2021. A unified analysis of a class of proximal bundle methods for solving hybrid convex composite optimization problems.

# Single Cut Model in the Stochastic Setting

Aggregate all cuts into a single one

$$\Gamma^+(u) = \tau\Gamma(u) + (1-\tau)[F(x,\xi) + \langle s(x,\xi), u - x\rangle].$$

Since

$$\mathbb{E}[F(x,\xi) + \langle s(x,\xi), u - x\rangle] = f(x) + \langle f'(x), u - x\rangle \leq f(u),$$

we have by induction

$$\mathbb{E}[\Gamma^+(u)] \leq f(u).$$

# Stochastic Composite Proximal Bundle (SCPB) Framework

1. Let $\lambda, \theta > 0$, integer $K \geq 1$, and $x_0 \in \operatorname{dom} h$ be given, and set $x_0^c = x_0$, $j = k = 1$, $j_0 = 0$, and

$$\tau = \frac{\theta K}{\theta K + 1};$$

2. Take an independent sample $\xi_{j-1}$ of r.v. $\xi$, set

$$x_j^c = \left\{ \begin{array}{ll} x_{j_{k-1}}, & \text{if } j = j_{k-1} + 1, \\ x_{j-1}^c, & \text{otherwise,} \end{array} \right.$$

and compute

$$x_j = \operatorname*{argmin}_{u \in \mathbb{R}^n} \left\{ h(u) + \langle S_j, u \rangle + \frac{1}{2\lambda} \|u - x_j^c\|^2 \right\},$$

where

$$S_j := \left\{ \begin{array}{ll} s(x_{j_{k-1}}, \xi_{j_{k-1}}), & \text{if } j = j_{k-1} + 1, \\ (1 - \tau)s(x_{j-1}, \xi_{j-1}) + \tau S_{j-1}, & \text{otherwise,} \end{array} \right.$$

3. Compute

$$y_j = \begin{cases} x_j, & \text{if } j = j_{k-1} + 1, \\ (1-\tau)x_j + \tau y_{j-1}, & \text{otherwise}; \end{cases}$$

4. Choose an integer $j_k \geq j_{k-1} + 1$, and set $\hat{y}_k = y_{j_k}$ when the $k$-th cycle ends;
5. if $k = K$ then **stop** and output

$$\hat{y}_K^a = \frac{1}{\lceil K/2 \rceil} \sum_{k=\lfloor K/2 \rfloor + 1}^{K} \hat{y}_k;$$

otherwise, set $k \leftarrow k + 1$ and $j \leftarrow j + 1$, and go to step 1.

# Remarks on SCPB

- An aggregated single cut
- No termination criterion for a cycle

Define a cycle
$$\mathcal{C}_k := \{i_k, \ldots, j_k\}, \quad \text{where} \quad i_k := j_{k-1} + 1$$

Two ways of setting $j_k$:

(B1) the smallest integer $j_k \geq i_k$ and $\lambda k \tau^{j_k - i_k} \leq C$;

(B2) the smallest integer $j_k \geq i_k + 1$ and $\lambda k \tau^{j_k - i_k} t_{i_k} \leq C$.

(B1) is deterministic and (B2) is stochastic

# Main Results – SCPB based on (B1)

Assume that conditions (A1)-(A4) hold and $\operatorname{dom} h$ has a finite diameter $D > 0$.

SCPB1 satisfies the following statements:

- Number of iterations within $\mathcal{C}_k$, or number of null steps

$$|\mathcal{C}_k| \leq \left\lceil (\theta K + 1) \ln \left( \frac{\lambda k}{C} + 1 \right) \right\rceil + 1.$$

- Convergence of SCPB1

$$\mathbb{E}[\phi(\hat{y}_K^a)] - \phi_* \leq \frac{1}{K} \left( \frac{D^2}{\lambda} + \frac{6C \min\{\lambda M^2, MD\}}{\lambda} + \frac{2\lambda M^2}{\theta} \right).$$

# A Practical Variant of SCPB1

Let pair $(\lambda, K)$ and constant $m \geq 1$ be given, and define

$$\theta = \frac{m}{K}, \quad C = \frac{D}{6M},$$

SCPB1 satisfies the following statements:

- Number of iterations within $\mathcal{C}_k$, or number of null steps

$$|\mathcal{C}_k| \leq \left\lceil (m+1) \ln \left( \frac{\lambda k}{C} + 1 \right) \right\rceil + 1.$$

- Convergence of SCPB1

$$\mathbb{E}[\phi(\hat{y}_K^a)] - \phi_* \leq \frac{2D^2}{\lambda K} + \frac{2\lambda M^2}{m}.$$

- Its expected overall iteration complexity is $\tilde{\mathcal{O}}(mK)$.

$$x_j = \operatorname*{argmin}_{u \in X} \left\{ \langle s(x_{j-1}, \xi_{j-1}), u \rangle + \frac{1}{2\lambda} \|u - x_{j-1}\|^2 \right\} \qquad \forall j = 1, \dots, K.$$

- Convergence of RSA

$$\mathbb{E}[\phi(x_K^a)] - \phi_* \leq \frac{2D^2}{\lambda K} + 2\lambda M^2, \quad x_K^a = \frac{1}{\lceil K/2 \rceil} \sum_{j=\lfloor K/2 \rfloor + 1}^{K} x_j.$$

Taking $\lambda = \frac{\sqrt{m}D}{M\sqrt{K}}$, given $\varepsilon > 0$, to obtain $x \in \operatorname{dom} h$ such that $\mathbb{E}[\phi(x)] - \phi_* \leq \varepsilon$,

- RSA has iteration complexity $\mathcal{O}\left(\frac{mM^2D^2}{\varepsilon^2}\right)$;
- SCPB1 has iteration complexity $\tilde{\mathcal{O}}\left(\frac{M^2D^2}{\varepsilon^2}\right)$.

---

[2] Nemirovski, Juditsky, Lan and Shapiro, 2009. Robust stochastic approximation approach to stochastic programming.

# Relationship between SCPB1 and RSA

Recall (B1) the smallest integer $j_k \geq i_k$ and $\lambda k \tau^{j_k - i_k} \leq C$.
Choosing

$$C = \frac{\alpha D \sqrt{K}}{M}, \quad \lambda = \frac{\alpha D}{M \sqrt{K}},$$

then (B1) is satisfied with $j_k = i_k$, since

$$\frac{C}{\lambda k} \geq \frac{C}{\lambda K} = 1 = \tau^{j_k - i_k}.$$

In summary,

- RSA performs one iteration per cycle
- RSA $\rightarrow$ SCPB1 is analogous to Subgradient method $\rightarrow$ PB
- RSA is restricted to small stepsizes, while SCPB1 can use large ones
- SCPB1 reduces the variance and the sample complexity by $m$

RSA: $\mathbb{E}[\phi(x_K^a)] - \phi_* \leq \frac{2D^2}{\lambda K} + 2\lambda M^2$, SCPB1: $\mathbb{E}[\phi(\hat{y}_K^a)] - \phi_* \leq \frac{2D^2}{\lambda K} + \frac{2\lambda M^2}{m}$

# Main Results – SCPB based on (B2)

Recall (B2): the smallest integer $j_k \geq i_k + 1$ and $\lambda k \tau^{j_k - i_k} t_{i_k} \leq C$.

Assume that conditions (A1)-(A4) hold and $\operatorname{dom} h$ has a finite diameter $D > 0$.

SCPB2 satisfies the following statements:

- Number of iterations within $\mathcal{C}_k$, or number of null steps

$$|\mathcal{C}_k| \leq \left\lceil (\theta K + 1) \ln \left( \frac{2M^2 \lambda^2 k}{C} + 1 \right) \right\rceil + 1.$$

- Convergence of SCPB2

$$\mathbb{E}[\phi(\hat{y}_K^a)] - \phi_* \leq \frac{1}{K} \left( \frac{3C + D^2}{\lambda} + \frac{2\lambda M^2}{\theta} + \frac{2\lambda M^2}{\theta^2 K} \right).$$

# A Practical Variant of SCPB2

Let pair $(\lambda, K)$ and constant $m \geq 1$ be given, and define

$$\theta = \frac{m}{K}, \quad C = \frac{D^2}{3},$$

SCPB2 satisfies the following statements:

- Number of iterations within $\mathcal{C}_k$, or number of null steps

$$|\mathcal{C}_k| \leq \left\lceil (m+1) \ln \left( \frac{6M^2 \lambda^2 k}{D^2} + 1 \right) \right\rceil + 1.$$

- Convergence of SCPB2

$$\mathbb{E}[\phi(\hat{y}_K^a)] - \phi_* \leq \frac{2D^2}{\lambda K} + \frac{4\lambda M^2}{m}.$$

- Its expected overall iteration complexity is $\tilde{\mathcal{O}}(mK)$.

# Test 1 – Two-stage Stochastic Program

$$\left\{ \begin{array}{l} \min \ c^T x_1 + \mathbb{E}[Q(x_1, \xi)] \\ x_1 \in \mathbb{R}^n : x_1 \geq 0, \sum_{i=1}^n x_1(i) = 1 \end{array} \right.$$

where the second stage recourse function is given by

$$Q(x_1, \xi) = \left\{ \begin{array}{l} \min\limits_{x_2 \in \mathbb{R}^n} \ \dfrac{1}{2} \left( \begin{array}{c} x_1 \\ x_2 \end{array} \right)^T \left( \xi \xi^T + \lambda_0 I_{2n} \right) \left( \begin{array}{c} x_1 \\ x_2 \end{array} \right) + \xi^T \left( \begin{array}{c} x_1 \\ x_2 \end{array} \right) \\ x_2 \geq 0, \sum_{i=1}^n x_2(i) = 1. \end{array} \right.$$

Table: $n = 50$, $N = 4000$

| Statistics | RSA | SCPB1 | SCPB2 |
|------------|-----|-------|-------|
| $\lambda$ | $7.4 \times 10^{-7}$ | $10^{-3}$ | $10^{-3}$ |
| Min Inner | 1 | 9 | 2 |
| Max Inner | 1 | 52 | 43 |
| Avg Inner | 1 | 43 | 5 |

Prob1 RSA vs SCPB1 vs SCPB2

Prob1 RSA vs SCPB1 vs SCPB2 new

# Test 2 – Two-stage Stochastic Program

$$\begin{cases} \min \ c^T x_1 + \mathbb{E}[\mathfrak{Q}(x_1, \xi)] \\ x_1 \in \mathbb{R}^n : \|x_1 - x_0\|_2 \leq 1 \end{cases}$$
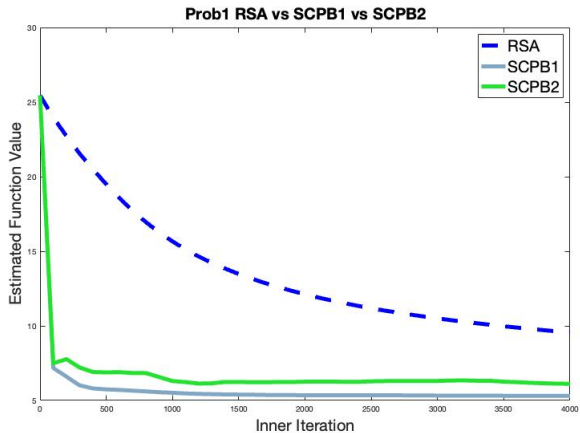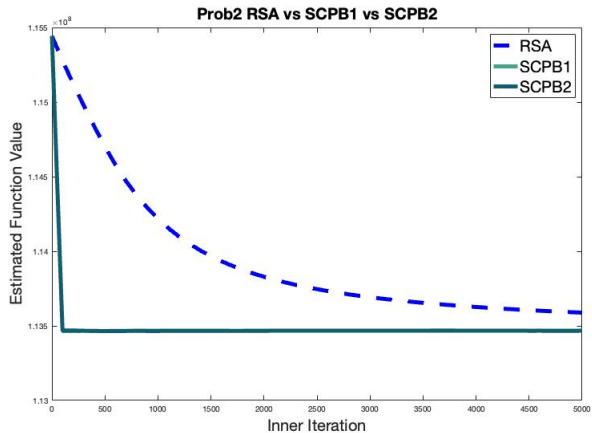
where the second stage recourse function is given by

$$Q(x_1, \xi) = \begin{cases} \min_{x_2 \in \mathbb{R}^n} \dfrac{1}{2} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}^T \left(\xi\xi^T + \lambda_0 I_{2n}\right) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \xi^T \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \\ \|x_2 - y_0\|_2^2 + \|x_1 - x_0\|_2^2 - R^2 \leq 0. \end{cases}$$

Table: $n = 50$, $N = 5000$

| Statistics | RSA | SCPB1 | SCPB2 |
|---|---|---|---|
| $\lambda$ | $8.9 \times 10^{-10}$ | $10^{-3}$ | $10^{-3}$ |
| Min Inner | 1 | 71 | 54 |
| Max Inner | 1 | 109 | 89 |
| Avg Inner | 1 | 100 | 77 |

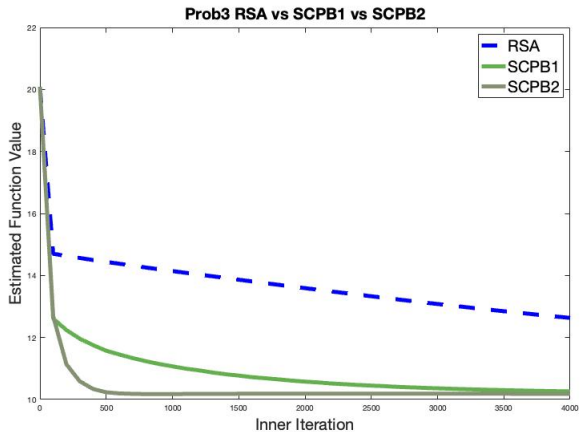Prob2 RSA vs SCPB1 vs SCPB2

# Test 3 – One-stage Stochastic Program

$$\min_{x \in X} \mathbb{E}\left[ \phi\left( \sum_{i=1}^{n} (\frac{i}{n} + \xi_i) x_i \right) \right]$$

where $X$ is the unit simplex.

Table: $n = 100$, $N = 4000$

| Statistics | RSA | SCPB1 | SCPB2 |
|------------|-----|-------|-------|
| $\lambda$ | $2.8 \times 10^{-5}$ | $10^{-3}$ | $10^{-3}$ |
| Min Inner | 1 | 1 | 2 |
| Max Inner | 1 | 26 | 6 |
| Avg Inner | 1 | 17 | 2 |

# Take-away Message

- The first proximal bundle method for stochastic programming

- A single cut aggregating all past information

- Optimal complexity for large stepsizes

- Includes RSA as an instance while outperforms RSA

- Variance reduction

# Reference

J. Liang, V. Guigues and R. D. C. Monteiro. A single cut proximal bundle method for stochastic convex composite optimization. ArXiv:2207.09024, 2022.

# Thank you!