

Algorithms for Deterministically Constrained Stochastic Optimization

Frank E. Curtis, Lehigh University

involving joint work with

Albert Berahas, University of Michigan

Michael O'Neill, Lehigh University (soon UNC Chapel Hill)

Daniel P. Robinson, Lehigh University

Baoyu Zhou, Lehigh University (soon Chicago Booth)

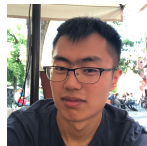
presented at

Center for Applied Mathematics, Cornell University

March 25, 2022



Collaborators and references



- ▶ A. S. Berahas, F. E. Curtis, D. P. Robinson, and B. Zhou, “Sequential Quadratic Optimization for Nonlinear Equality Constrained Stochastic Optimization,” *SIAM Journal on Optimization*, 31(2):1352–1379, 2021.
- ▶ A. S. Berahas, F. E. Curtis, M. J. O’Neill, and D. P. Robinson, “A Stochastic Sequential Quadratic Optimization Algorithm for Nonlinear Equality Constrained Optimization with Rank-Deficient Jacobians,” <https://arxiv.org/abs/2106.13015>.
- ▶ F. E. Curtis, D. P. Robinson, and B. Zhou, “Inexact Sequential Quadratic Optimization for Minimizing a Stochastic Objective Subject to Deterministic Nonlinear Equality Constraints,” <https://arxiv.org/abs/2107.03512>.
- ▶ F. E. Curtis, M. J. O’Neill, and D. P. Robinson, “Worst-Case Complexity of an SQP Method for Nonlinear Equality Constrained Stochastic Optimization,” <https://arxiv.org/abs/2112.14799>.

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

Worst-Case Iteration Complexity

Extensions

Conclusion

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

Worst-Case Iteration Complexity

Extensions

Conclusion

Constrained optimization (deterministic)

Consider

$$\begin{array}{ll}\min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & c_{\mathcal{E}}(x) = 0 \\ & c_{\mathcal{I}}(x) \leq 0\end{array}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $c_{\mathcal{E}} : \mathbb{R}^n \rightarrow \mathbb{R}^{m_{\mathcal{E}}}$, and $c_{\mathcal{I}} : \mathbb{R}^n \rightarrow \mathbb{R}^{m_{\mathcal{I}}}$ are smooth

- ▶ Physics-constrained, resource-constrained, etc.
- ▶ Long history of algorithms (penalty, SQP, interior-point, etc.)
- ▶ Comprehensive theory (even with lack of constraint qualifications)
- ▶ Effective software (Ipopt, Knitro, LOQO, etc.)

Constrained optimization (stochastic constraints)

Consider

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & c_{\mathcal{E}}(x) = 0 \\ & c_{\mathcal{I}}(x, \omega) \lesssim 0 \end{array}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $c_{\mathcal{E}} : \mathbb{R}^n \rightarrow \mathbb{R}^{m_{\mathcal{E}}}$, and $c_{\mathcal{I}} : \mathbb{R}^n \times \Omega \rightarrow \mathbb{R}^{m_{\mathcal{I}}}$

- ▶ Various modeling paradigms:
- ▶ ...stochastic optimization
- ▶ ... (distributionally) robust optimization
- ▶ ...chance-constrained optimization

Constrained optimization (stochastic objective)

Consider

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & f(x) \equiv \mathbb{E}[F(x, \omega)] \\ \text{s.t.} & c_{\mathcal{E}}(x) = 0 \\ & c_{\mathcal{I}}(x) \leq 0 \end{array}$$

where $f : \mathbb{R}^n \times \mathbb{R}$, $F : \mathbb{R}^n \times \Omega \rightarrow \mathbb{R}$, $c_{\mathcal{E}} : \mathbb{R}^n \rightarrow \mathbb{R}^{m_{\mathcal{E}}}$, and $c_{\mathcal{I}} : \mathbb{R}^n \rightarrow \mathbb{R}^{m_{\mathcal{I}}}$

- ▶ ω has probability space (Ω, \mathcal{F}, P)
- ▶ $\mathbb{E}[\cdot]$ with respect to P
- ▶ Classical applications under uncertainty, constrained DNN training, etc.
- ▶ Besides cases involving a deterministic equivalent...
- ▶ ... very few algorithms so far (mostly penalty methods)

What kind of algorithm do we want?

Need to establish what we want/expect from an algorithm.

Note: We are interested in the **fully stochastic** regime.[†]

We assume:

- ▶ Feasible methods are not tractable
- ▶ ... so no projection methods, Frank-Wolfe, etc.
- ▶ “Two-phase” methods are not effective
- ▶ ... so should not search for feasibility, then optimize.
- ▶ Only enforce convergence in expectation.

Finally, want to use techniques that can generalize to diverse settings.

[†] Alternatively, see Na, Animescu, Kolar (2021, 2022)

This talk

Consider *equality constrained* stochastic optimization:

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & f(x) \equiv \mathbb{E}[F(x, \omega)] \\ \text{s.t.} & c_{\mathcal{E}}(x) = 0 \end{array}$$

- ▶ *Adaptive SQP* method for deterministic setting
- ▶ *Stochastic SQP* method for stochastic setting
- ▶ Convergence in expectation (comparable to SG for unconstrained setting)
- ▶ Worst-case complexity on par with stochastic subgradient method
- ▶ Numerical experiments are very promising
- ▶ Various open questions!

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

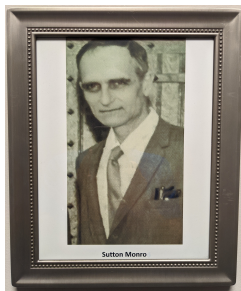
Worst-Case Iteration Complexity

Extensions

Conclusion

Stochastic gradient method (SG)

Invented by Herbert Robbins and Sutton Monro (1951)



Sutton Monro, former Lehigh faculty member

Stochastic gradient (*not* descent)

$$\min_{x \in \mathbb{R}^n} f(x) \equiv \mathbb{E}[F(x, \omega)]$$

where $g := \nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant L

Algorithm SG : Stochastic Gradient

- 1: choose an initial point $x_0 \in \mathbb{R}^n$ and step sizes $\{\alpha_k\} > 0$
 - 2: **for** $k \in \{0, 1, 2, \dots\}$ **do**
 - 3: set $x_{k+1} \leftarrow x_k - \alpha_k \bar{g}_k$, where $\mathbb{E}_k[\bar{g}_k] = g_k$ and $\mathbb{E}_k[\|\bar{g}_k - g_k\|_2^2] \leq M$
 - 4: **end for**
-

Not a descent method! ...but *eventual descent in expectation*:

$$\begin{aligned} f(x_{k+1}) - f(x_k) &\leq g_k^T(x_{k+1} - x_k) + \frac{1}{2}L\|x_{k+1} - x_k\|_2^2 \\ &= -\alpha_k g_k^T \bar{g}_k + \frac{1}{2}\alpha_k^2 L\|\bar{g}_k\|_2^2 \\ \implies \mathbb{E}_k[f(x_{k+1})] - f(x_k) &\leq -\alpha_k \|g_k\|_2^2 + \frac{1}{2}\alpha_k^2 L\mathbb{E}_k[\|\bar{g}_k\|_2^2]. \end{aligned}$$

Markovian: x_{k+1} depends only on x_k and random choice at iteration k .

SG theory

Theorem SG

Since $\mathbb{E}_k[\bar{g}_k] = g_k$ and $\mathbb{E}_k[\|\bar{g}_k - g_k\|_2^2] \leq M$ for all $k \in \mathbb{N}$:

$$\alpha_k = \frac{1}{L} \quad \Rightarrow \quad \mathbb{E} \left[\frac{1}{k} \sum_{j=1}^k \|g_j\|_2^2 \right] \leq \mathcal{O}(M)$$

$$\alpha_k = \Theta \left(\frac{1}{k} \right) \quad \Rightarrow \quad \mathbb{E} \left[\frac{1}{\left(\sum_{j=1}^k \alpha_j \right)} \sum_{j=1}^k \alpha_j \|g_j\|_2^2 \right] \rightarrow 0$$

$$\Rightarrow \liminf_{k \rightarrow \infty} \mathbb{E}[\|g_k\|_2^2] = 0$$

SG illustration

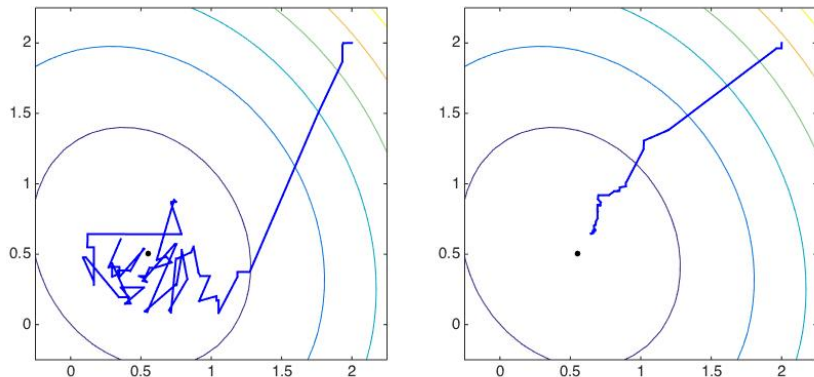


Figure: SG with fixed step size (left) vs. diminishing step sizes (right)

Sequential quadratic optimization (SQP)

Consider

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & c(x) = 0 \end{array}$$

with $g \equiv \nabla f$, $J \equiv \nabla c$, and $H \succ 0$ (for simplicity), two viewpoints:

$$\begin{bmatrix} g(x) + J(x)^T y \\ c(x) \end{bmatrix} = 0$$

or

$$\begin{array}{ll} \min_{d \in \mathbb{R}^n} & f(x) + g(x)^T d + \frac{1}{2} d^T H d \\ \text{s.t.} & c(x) + J(x)d = 0 \end{array}$$

both leading to the same “Newton-SQP system”:

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} d_k \\ y_k \end{bmatrix} = - \begin{bmatrix} g_k \\ c_k \end{bmatrix}$$

SQP illustration

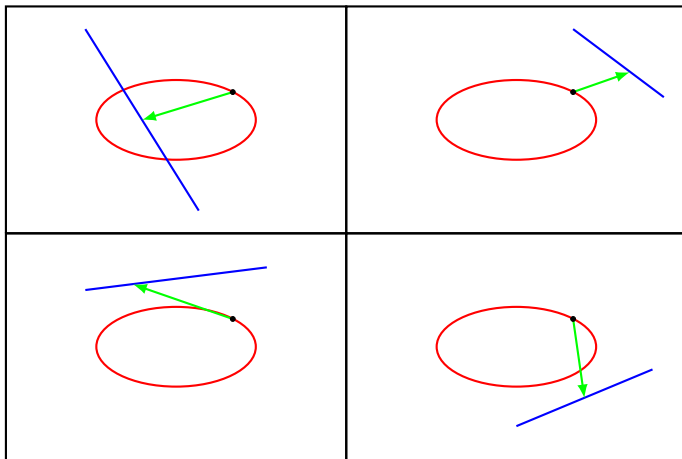


Figure: Illustrations of SQP subproblem solutions

SQP

- Algorithm guided by merit function, with **adaptive** parameter τ , defined by

$$\phi(x, \tau) = \tau f(x) + \|c(x)\|_1$$

a model of which is defined as

$$q(x, \tau, d) = \tau(f(x) + g(x)^T d + \frac{1}{2} d^T H d) + \|c(x) + J(x)d\|_1$$

- For a given $d \in \mathbb{R}^n$ satisfying $c(x) + J(x)d = 0$, the reduction in this model is

$$\Delta q(x, \tau, d) = -\tau(g(x)^T d + \frac{1}{2} d^T H d) + \|c(x)\|_1,$$

and it is easily shown that

$$\phi'(x, \tau, d) \leq -\Delta q(x, \tau, d)$$

SQP with backtracking line search

Algorithm SQP-B

1: choose $x_0 \in \mathbb{R}^n$, $\tau_{-1} \in \mathbb{R}_{>0}$, $\sigma \in (0, 1)$, $\eta \in (0, 1)$ 2: **for** $k \in \{0, 1, 2, \dots\}$ **do**

3: solve

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} d_k \\ y_k \end{bmatrix} = - \begin{bmatrix} g_k \\ c_k \end{bmatrix}$$

4: set τ_k to ensure $\Delta q(x_k, \tau_k, d_k) \gg 0$, offered by

$$\tau_k \leq \frac{(1 - \sigma)\|c_k\|_1}{g_k^T d_k + d_k^T H_k d_k} \quad \text{if } g_k^T d_k + d_k^T H_k d_k > 0$$

5: backtracking line search to ensure $x_{k+1} \leftarrow x_k + \alpha_k d_k$ yields

$$\phi(x_{k+1}, \tau_k) \leq \phi(x_k, \tau_k) - \eta \alpha_k \Delta q(x_k, \tau_k, d_k)$$

6: **end for**

Convergence theory

Assumption

- ▶ f , c , g , and J bounded and Lipschitz
- ▶ singular values of J bounded below (i.e., the LICQ)
- ▶ $u^T H_k u \geq \zeta \|u\|_2^2$ for all $u \in \text{Null}(J_k)$ for all $k \in \mathbb{N}$

Theorem SQP-B

- ▶ $\{\alpha_k\} \geq \alpha_{\min}$ for some $\alpha_{\min} > 0$
- ▶ $\{\tau_k\} \geq \tau_{\min}$ for some $\tau_{\min} > 0$
- ▶ $\Delta q(x_k, \tau_k, d_k) \rightarrow 0$ implies

$$\|d_k\|_2 \rightarrow 0, \quad \|c_k\|_2 \rightarrow 0, \quad \|g_k + J_k^T y_k\|_2 \rightarrow 0$$

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

Worst-Case Iteration Complexity

Extensions

Conclusion

Toward stochastic SQP

- ▶ In a stochastic setting, line searches are (likely) intractable
- ▶ However, for ∇f and ∇c , may have Lipschitz constants (or estimates)
- ▶ Step #1: Design an **adaptive** SQP method with

step sizes determined by Lipschitz constant estimates

- ▶ Step #2: Design a **stochastic** SQP method on this approach

Primary challenge: Nonsmoothness

In SQP-B, step size is chosen based on reducing the merit function.

The merit function is nonsmooth! An upper bound is

$$\begin{aligned} & \phi(x_k + \alpha_k d_k, \tau_k) - \phi(x_k, \tau_k) \\ & \leq \alpha_k \tau_k g_k^T d_k + |1 - \alpha_k| \|c_k\|_1 - \|c_k\|_1 + \frac{1}{2}(\tau_k L_k + \Gamma_k) \alpha_k^2 \|d_k\|_2^2 \end{aligned}$$

where L_k and Γ_k are Lipschitz constant estimates for f and $\|c\|_1$ at x_k

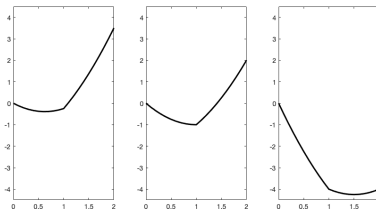


Figure: Three cases for upper bound of ϕ

Idea: Choose α_k to ensure sufficient decrease using this bound

SQP with adaptive step sizes

Algorithm SQP-A

1: choose $x_0 \in \mathbb{R}^n$, $\tau_{-1} \in \mathbb{R}_{>0}$, $\sigma \in (0, 1)$, $\eta \in (0, 1)$

2: **for** $k \in \{0, 1, 2, \dots\}$ **do**

3: solve

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} d_k \\ y_k \end{bmatrix} = - \begin{bmatrix} g_k \\ c_k \end{bmatrix}$$

4: set τ_k to ensure $\Delta q(x_k, \tau_k, d_k) \gg 0$, offered by

$$\tau_k \leq \frac{(1 - \sigma) \|c_k\|_1}{g_k^T d_k + d_k^T H_k d_k} \quad \text{if } g_k^T d_k + d_k^T H_k d_k > 0$$

5: set

$$\hat{\alpha}_k \leftarrow \frac{2(1 - \eta) \Delta q(x_k, \tau_k, d_k)}{(\tau_k L_k + \Gamma_k) \|d_k\|_2^2} \quad \text{and}$$

$$\tilde{\alpha}_k \leftarrow \hat{\alpha}_k - \frac{4 \|c_k\|_1}{(\tau_k L_k + \Gamma_k) \|d_k\|_2^2}$$

6: set

$$\alpha_k \leftarrow \begin{cases} \hat{\alpha}_k & \text{if } \hat{\alpha}_k < 1 \\ 1 & \text{if } \tilde{\alpha}_k \leq 1 \leq \hat{\alpha}_k \\ \tilde{\alpha}_k & \text{if } \tilde{\alpha}_k > 1 \end{cases}$$

7: set $x_{k+1} \leftarrow x_k + \alpha_k d_k$ and continue or update L_k and/or Γ_k and return to step 5

8: **end for**

Convergence theory

Exactly the same as for SQP-B, except different step size lower bound

- For SQP-A:

$$\alpha_k = \frac{2(1-\eta)\Delta q(x_k, \tau_k, d_k)}{(\tau_k L_k + \Gamma_k)\|d_k\|_2^2} \geq \frac{2(1-\eta)\kappa_q \tau_{\min}}{(\tau_{-1}\rho L + \rho\Gamma)\kappa_\Psi} > 0$$

- For SQP-B:

$$\alpha_k > \frac{2\nu(1-\eta)\Delta q(x_k, \tau_k, d_k)}{(\tau_k L + \Gamma)\|d_k\|_2^2} \geq \frac{2\nu(1-\eta)\kappa_q \tau_{\min}}{(\tau_{-1}L + \Gamma)\kappa_\Psi} > 0$$

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

Worst-Case Iteration Complexity

Extensions

Conclusion

Stochastic setting

Consider the stochastic problem:

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & f(x) \equiv \mathbb{E}[F(x, \omega)] \\ \text{s.t.} & c(x) = 0 \end{array}$$

Let us assume only the following:

Assumption

For all $k \in \mathbb{N}$, one can compute \bar{g}_k with

$$\mathbb{E}_k[\bar{g}_k] = g_k \quad \text{and} \quad \mathbb{E}_k[\|\bar{g}_k - g_k\|_2^2] \leq M$$

Search directions computed by:

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} \bar{d}_k \\ \bar{y}_k \end{bmatrix} = - \begin{bmatrix} \bar{g}_k \\ c_k \end{bmatrix}$$

Important: Given x_k , the values (c_k, J_k, H_k) are **determined**

Stochastic SQP with adaptive step sizes

(For simplicity, assume Lipschitz constants L and Γ are known.)

Algorithm : Stochastic SQP

1: choose $x_0 \in \mathbb{R}^n$, $\bar{\tau}_{-1} \in \mathbb{R}_{>0}$, $\sigma \in (0, 1)$, $\{\beta_k\} \in (0, 1)$

2: **for** $k \in \{0, 1, 2, \dots\}$ **do**

3: solve

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} \bar{d}_k \\ \bar{y}_k \end{bmatrix} = - \begin{bmatrix} \bar{g}_k \\ c_k \end{bmatrix}$$

4: set $\bar{\tau}_k$ to ensure $\Delta \bar{q}(x_k, \bar{\tau}_k, \bar{d}_k) \gg 0$, offered by

$$\bar{\tau}_k \leq \frac{(1 - \sigma) \|c_k\|_1}{\bar{g}_k^T \bar{d}_k + \bar{d}_k^T H_k \bar{d}_k} \quad \text{if } \bar{g}_k^T \bar{d}_k + \bar{d}_k^T H_k \bar{d}_k > 0$$

5: set

$$\bar{\bar{\alpha}}_k \leftarrow \frac{\beta_k \Delta \bar{q}(x_k, \bar{\tau}_k, \bar{d}_k)}{(\bar{\tau}_k L + \Gamma) \|\bar{d}_k\|_2^2} \quad \text{and}$$

$$\bar{\bar{\alpha}}_k \leftarrow \bar{\bar{\alpha}}_k - \frac{4 \|c_k\|_1}{(\bar{\tau}_k L + \Gamma) \|\bar{d}_k\|_2^2}$$

6: set

$$\bar{\alpha}_k \leftarrow \begin{cases} \bar{\bar{\alpha}}_k & \text{if } \bar{\bar{\alpha}}_k < 1 \\ 1 & \text{if } \bar{\bar{\alpha}}_k \leq 1 \leq \bar{\bar{\alpha}}_k \\ \bar{\bar{\alpha}}_k & \text{if } \bar{\bar{\alpha}}_k > 1 \end{cases}$$

7: set $x_{k+1} \leftarrow x_k + \bar{\alpha}_k \bar{d}_k$

8: **end for**

step size control

The sequence $\{\beta_k\}$ allows us to consider, like for SG,

- ▶ a fixed step size
- ▶ diminishing step sizes (e.g., $\Theta(1/k)$)

Unfortunately, additional control on the step size is needed

- ▶ too small: insufficient progress
- ▶ too large: ruins progress toward feasibility / optimality

We never know when the step size is too small or too large!

Idea: Project $\tilde{\alpha}_k$ and $\bar{\alpha}_k$ onto

$$\left[\frac{\beta_k \bar{\tau}_k}{\bar{\tau}_k L + \Gamma}, \frac{\beta_k \bar{\tau}_k}{\bar{\tau}_k L + \Gamma} + \theta \beta_k^2 \right]$$

where $\theta \in \mathbb{R}_{>0}$ is a user-defined parameter

Fundamental lemmas

Lemma

For all $k \in \mathbb{N}$, for any realization of \bar{g}_k , one finds

$$\begin{aligned} & \phi(x_k + \bar{\alpha}_k \bar{d}_k, \bar{\tau}_k) - \phi(x_k, \bar{\tau}_k) \\ \leq & \underbrace{-\bar{\alpha}_k \Delta q(x_k, \bar{\tau}_k, d_k)}_{\mathcal{O}(\beta_k), \text{“deterministic”}} + \underbrace{\frac{1}{2} \bar{\alpha}_k \beta_k \Delta \bar{q}(x_k, \bar{\tau}_k, \bar{d}_k)}_{\mathcal{O}(\beta_k^2), \text{stochastic/noise}} + \underbrace{\bar{\alpha}_k \bar{\tau}_k g_k^T (\bar{d}_k - d_k)}_{\text{due to adaptive } \bar{\alpha}_k} \end{aligned}$$

Lemma

For all $k \in \mathbb{N}$, one finds

$$\mathbb{E}_k[\bar{d}_k] = d_k, \quad \mathbb{E}_k[\bar{y}_k] = y_k, \quad \text{and} \quad \mathbb{E}_k[\|\bar{d}_k - d_k\|_2] = \mathcal{O}(\sqrt{M})$$

as well as

$$g_k^T d_k \geq \mathbb{E}_k[\bar{g}_k^T \bar{d}_k] \geq g_k^T d_k - \zeta^{-1} M \quad \text{and} \quad d_k^T H_k d_k \leq \mathbb{E}_k[\bar{d}_k^T H_k \bar{d}_k]$$

Good merit parameter behavior

Lemma

If $\{\bar{\tau}_k\}$ eventually remains fixed at sufficiently small $\tau_{\min} > 0$, then for large k

$$\mathbb{E}_k[\bar{\alpha}_k \bar{\tau}_k g_k^T(\bar{d}_k - d_k)] = \beta_k^2 \tau_{\min} \mathcal{O}(\sqrt{M})$$

Theorem

If $\{\bar{\tau}_k\}$ eventually remains fixed at sufficiently small $\tau_{\min} > 0$, then for large k

$$\beta_k = \Theta(1) \implies \alpha_k = \frac{\tau_{\min}}{\tau_{\min} L + \Gamma} \implies \mathbb{E} \left[\frac{1}{k} \sum_{j=1}^k \Delta q(x_j, \tau_{\min}, d_j) \right] \leq \mathcal{O}(M)$$

$$\beta_k = \Theta\left(\frac{1}{k}\right) \implies \mathbb{E} \left[\frac{1}{\left(\sum_{j=1}^k \beta_j\right)} \sum_{j=1}^k \beta_j \Delta q(x_j, \tau_{\min}, d_j) \right] \rightarrow 0$$

Good merit parameter behavior

Lemma

If $\{\bar{\tau}_k\}$ eventually remains fixed at sufficiently small $\tau_{\min} > 0$, then for large k

$$\mathbb{E}_k[\bar{\alpha}_k \bar{\tau}_k g_k^T (\bar{d}_k - d_k)] = \beta_k^2 \tau_{\min} \mathcal{O}(\sqrt{M})$$

Theorem

If $\{\bar{\tau}_k\}$ eventually remains fixed at sufficiently small $\tau_{\min} > 0$, then for large k

$$\beta_k = \Theta(1) \implies \alpha_k = \frac{\tau_{\min}}{\tau_{\min} L + \Gamma} \implies \mathbb{E} \left[\frac{1}{k} \sum_{j=1}^k (\|g_j + J_j^T y_j\|_2 + \|c_j\|_2) \right] \leq \mathcal{O}(M)$$

$$\beta_k = \Theta\left(\frac{1}{k}\right) \implies \mathbb{E} \left[\frac{1}{\left(\sum_{j=1}^k \beta_j\right)} \sum_{j=1}^k \beta_j (\|g_j + J_j^T y_j\|_2 + \|c_j\|_2) \right] \rightarrow 0$$

Poor merit parameter behavior

$\{\bar{\tau}_k\} \searrow 0$:

- ▶ cannot occur if $\|\bar{g}_k - g_k\|_2$ is bounded uniformly
- ▶ occurs with small probability if distribution of \bar{g}_k has fast decay

$\{\bar{\tau}_k\}$ remains too large:

- ▶ can only occur if realization of $\{\bar{g}_k\}$ is *one-sided for all* k
- ▶ if there exists $p \in (0, 1]$ such that, for all k in infinite \mathcal{K} ,

$$\mathbb{P}_k \left[\bar{g}_k^T \bar{d}_k + \max\{\bar{d}_k^T H_k \bar{d}_k, 0\} \geq g_k^T d_k + \max\{d_k^T H_k d_k, 0\} \right] \geq p$$

then occurs with probability zero

Numerical results

Matlab software: <https://github.com/frankecurtis/StochasticSQP>

CUTE problems with noise added to gradients with different noise levels

- ▶ Stochastic SQP: 10^3 iterations
- ▶ Stochastic Subgradient: 10^4 iterations and tuned over 11 values of τ

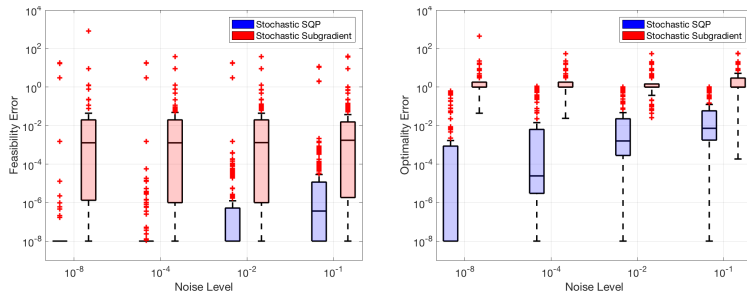


Figure: Box plots for feasibility errors (left) and optimality errors (right).

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

Worst-Case Iteration Complexity

Extensions

Conclusion

Complexity of deterministic algorithm

All reductions in the merit function can be cast in terms of smallest τ .

Lemma 7

If $\{\tau_k\}$ eventually remains fixed at sufficiently small τ_{\min} , then for any $\epsilon \in (0, 1)$ there exists $(\kappa_1, \kappa_2) \in (0, \infty) \times (0, \infty)$ such that, for all k ,

$$\|g_k + J_k^T y_k\| > \epsilon \text{ or } \sqrt{\|c_k\|_1} > \epsilon \implies \Delta q(x_k, \tau_k, d_k) \geq \min\{\kappa_1, \kappa_2 \tau_{\min}\} \epsilon.$$

Since τ_{\min} is determined by the initial point, *it will be reached.*

Theorem 8

For any $\epsilon \in (0, 1)$, there exists $(\kappa_1, \kappa_2) \in (0, \infty) \times (0, \infty)$ such that

$$\|g_k + J_k^T y_k\| \leq \epsilon \text{ and } \sqrt{\|c_k\|_1} \leq \epsilon$$

in a number of iterations no more than

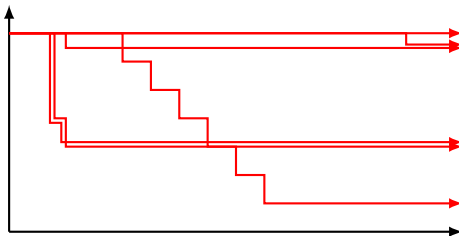
$$\left(\frac{\tau_{-1}(f_0 - f_{\inf}) + \|c_0\|_1}{\min\{\kappa_1, \kappa_2 \tau_{\min}\}} \right) \epsilon^{-2}.$$

Challenge in the stochastic setting

We are minimizing a function that is changing during the optimization.

In the stochastic setting, minimum τ is not determined by the initial point.

- ▶ Even if we assume $\tau_k \geq \tau_{\min} > 0$ for all k in any realization, the final value of the merit parameter τ is not determined.
- ▶ This means we cannot cast all reductions in terms of some fixed τ .



Our approach

In fact, τ reaching some minimum value is not necessary.

- ▶ Important: Diminishing probability of continued imbalance between “true” merit parameter update and “stochastic” merit parameter update.
- ▶ In iteration k , the algorithm has obtained the merit parameter value $\bar{\tau}_{k-1}$.
- ▶ If the true gradient is computed, then one obtains $\tau_k^{\text{trial, true}}$.
- ▶ Consider the random index set

$$\mathcal{K}_\tau := \{k : \tau_k^{\text{trial, true}} < \bar{\tau}_{k-1}\}.$$

Lemma 9

For any $\delta \in (0, 1)$, one finds that

$$\mathbb{P} \left[|\mathcal{K}_\tau| \leq \left\lceil \frac{\ell(s_{\max}, \delta)}{p} \right\rceil \right] \geq 1 - \delta,$$

where

$$\ell(s, \delta) := s + \log(1/\delta) + \sqrt{\log(1/\delta)^2 + 2s \log(1/\delta)} > 0.$$

Chernoff bound

How do we get there?

Lemma 10 (Chernoff bound, multiplicative form)

For any k , let $\{Y_0, \dots, Y_k\}$ be independent Bernoulli random variables. Then, for any $s_{\max} \in \mathbb{N}$ and $\delta \in (0, 1)$,

$$\sum_{j=0}^k \mathbb{P}[Y_j = 1] \geq \ell(s_{\max}, \delta) \implies \mathbb{P} \left[\sum_{j=0}^k Y_j \leq s_{\max} \right] \leq \delta.$$

We construct a tree whose nodes are signatures of possible runs of the algorithm.

- ▶ A realization $\{\bar{g}_0, \dots, \bar{g}_k\}$ belongs to a node if and only if a certain number of decreases of τ have occurred and the probability of decrease in the current iteration is in a given closed/open interval.
- ▶ Bad leaves are those when the probability of decrease has accumulated beyond a threshold, yet the merit parameter has not been decrease sufficiently often.
- ▶ Along the way, we apply a Chernoff bound on a carefully constructed set of random variables to bound probabilities associated with bad leaves.

Worst-case iteration complexity of $\tilde{\mathcal{O}}(\epsilon^{-4})$

Theorem 11

Suppose the algorithm is run

- ▶ k_{\max} iterations with
- ▶ $\beta_k = \gamma/\sqrt{k_{\max} + 1}$ and
- ▶ the merit parameter is reduced at most $s_{\max} \in \{0, 1, \dots, k_{\max}\}$ times.

Let k_* be sampled uniformly over $\{1, \dots, k_{\max}\}$. Then, with probability $1 - \delta$,

$$\begin{aligned} \mathbb{E}[\|g_{k_*} + J_{k_*}^T y_{k_*}\|_2^2 + \|c_{k_*}\|_1] &\leq \frac{\tau_{-1}(f_0 - f_{\inf}) + \|c_0\|_1 + M}{\sqrt{k_{\max} + 1}} \\ &\quad + \frac{(\tau_{-1} - \tau_{\min})(s_{\max} \log(k_{\max}) + \log(1/\delta))}{\sqrt{k_{\max} + 1}} \end{aligned}$$

Theorem 12

If the stochastic gradient estimates are sub-Gaussian, then w.p. $1 - \bar{\delta}$

$$s_{\max} = \mathcal{O} \left(\log \left(\log \left(\frac{k_{\max}}{\bar{\delta}} \right) \right) \right).$$

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

Worst-Case Iteration Complexity

Extensions

Conclusion

Recent work (under review): No LICQ

Remove constraint qualification

- ▶ infeasible and/or degenerate problems
- ▶ step decomposition method

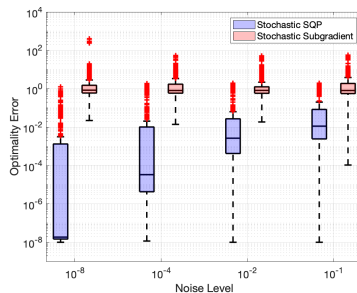
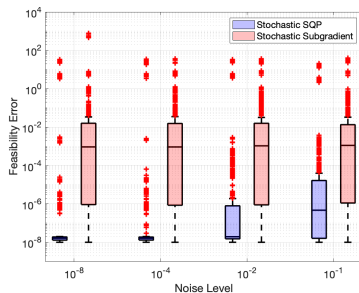


Figure: Box plots for feasibility errors (left) and optimality errors (right).

Recent work (under review): Matrix-free methods

Inexact subproblem solves

- ▶ stochasticity and inexactness(!)
- ▶ applicable for large-scale, e.g., PDE-constrained

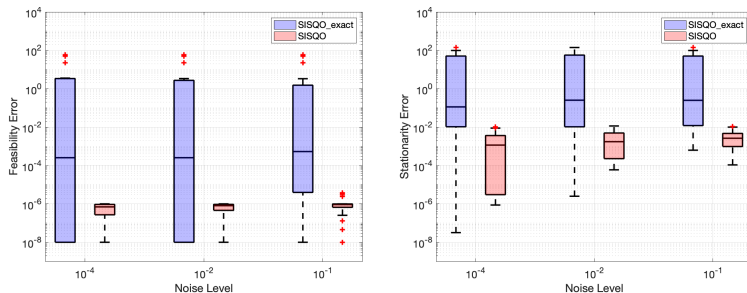


Figure: Box plots for feasibility errors (left) and optimality errors (right).

Current work: Inequality constraints

Inequality constraints

- ▶ SQP
- ▶ interior-point

Main challenge: For *equality* constraints only, subproblem solution on linearized constraints remains unbiased:

$$\begin{aligned} c_k + J_k \bar{d}_k = 0 & \iff \bar{d}_k = v_k + \bar{u}_k \\ & \text{with } v_k \in \text{Range}(J_k^T) \text{ and } \bar{u}_k \in \text{Null}(J_k) \\ & \text{has } E_k[\bar{u}_k] = u_k. \end{aligned}$$

However, when *inequalities* are present, subproblem solution is biased.

Outline

Motivation

SG and SQP

Adaptive (Deterministic) SQP

Stochastic SQP

Worst-Case Iteration Complexity

Extensions

Conclusion

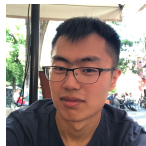
Summary

Consider *equality constrained* stochastic optimization:

$$\begin{array}{ll} \min_{x \in \mathbb{R}^n} & f(x) \equiv \mathbb{E}[F(x, \omega)] \\ \text{s.t.} & c_{\mathcal{E}}(x) = 0 \end{array}$$

- ▶ *Adaptive* SQP method for deterministic setting
- ▶ *Stochastic* SQP method for stochastic setting
- ▶ Convergence in expectation (comparable to SG for unconstrained setting)
- ▶ Worst-case complexity on par with stochastic subgradient method
- ▶ Numerical experiments are very promising
- ▶ Various extensions (on-going)

Collaborators and references



- ▶ A. S. Berahas, F. E. Curtis, D. P. Robinson, and B. Zhou, “Sequential Quadratic Optimization for Nonlinear Equality Constrained Stochastic Optimization,” *SIAM Journal on Optimization*, 31(2):1352–1379, 2021.
- ▶ A. S. Berahas, F. E. Curtis, M. J. O’Neill, and D. P. Robinson, “A Stochastic Sequential Quadratic Optimization Algorithm for Nonlinear Equality Constrained Optimization with Rank-Deficient Jacobians,” <https://arxiv.org/abs/2106.13015>.
- ▶ F. E. Curtis, D. P. Robinson, and B. Zhou, “Inexact Sequential Quadratic Optimization for Minimizing a Stochastic Objective Subject to Deterministic Nonlinear Equality Constraints,” <https://arxiv.org/abs/2107.03512>.
- ▶ F. E. Curtis, M. J. O’Neill, and D. P. Robinson, “Worst-Case Complexity of an SQP Method for Nonlinear Equality Constrained Stochastic Optimization,” <https://arxiv.org/abs/2112.14799>.

ICCOPT 2022

International Conference on Continuous Optimization



July 23–28, 2022

