

Algorithms for Deterministically Constrained Stochastic Optimization

Frank E. Curtis, Lehigh University

involving joint work with

Albert Berahas, University of Michigan

Michael O'Neill, UNC Chapel Hill

Daniel P. Robinson, Lehigh University

Baoyu Zhou, Chicago Booth

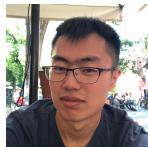
presented at

AIRS in the AIR

October 10, 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)

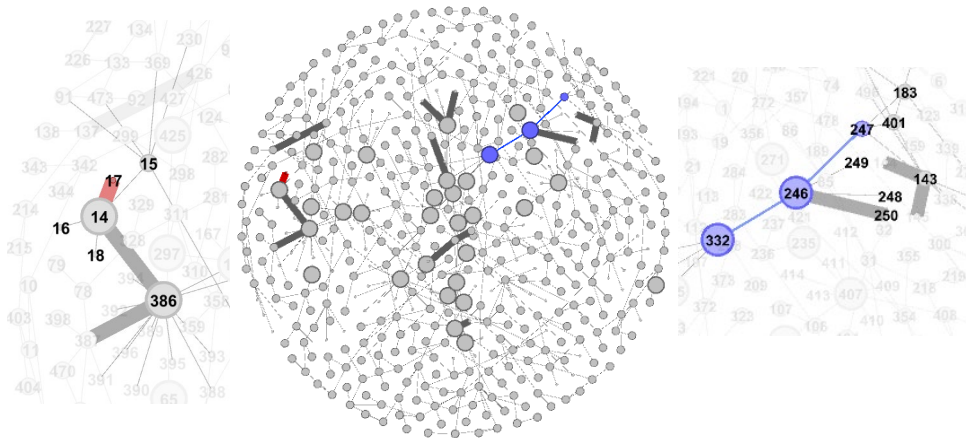
Consider

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

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

Motivation #1: Network optimization



Motivation #2: Physics-constrained learning

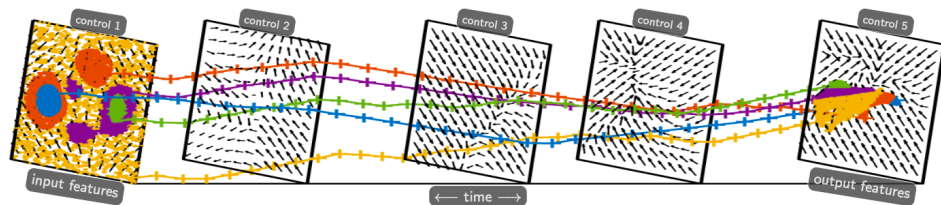
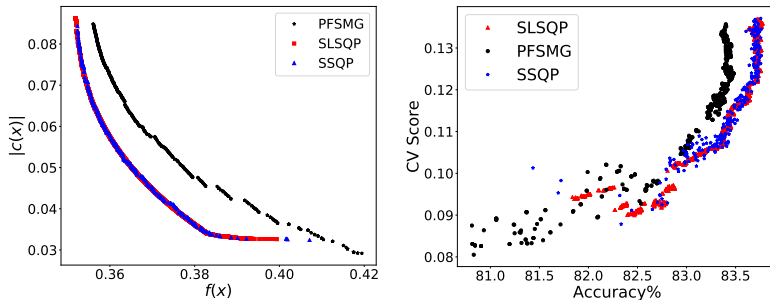


Photo: Lars Ruthotto, “An Optimal Control Framework for Efficient Training of Deep Neural Networks”

Motivation #3: Fair learning

$$\mathbb{P}[Y = y|A = 1] = \mathbb{P}[Y = y|A = 0] \text{ for each } y \in \{-1, 1\}$$



$$\min_{x \in \mathbb{R}^n} \frac{1}{N^o} \sum_{(v_i, y_i) \in D_o} \ell(x, v_i, y_i) \quad \text{s.t.} \quad -\epsilon \leq \frac{1}{N^c} \sum_{(v_i, a_i) \in D_c} (a_i - \bar{a}) x^T v_i \leq \epsilon$$

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, Anitescu, 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

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

$$\begin{bmatrix} \nabla f(x) + J(x)^T y \\ c(x) \end{bmatrix} = 0 \quad \text{or}$$

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

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

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} d_k \\ y_k \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) \\ c_k \end{bmatrix}$$

SQP illustration

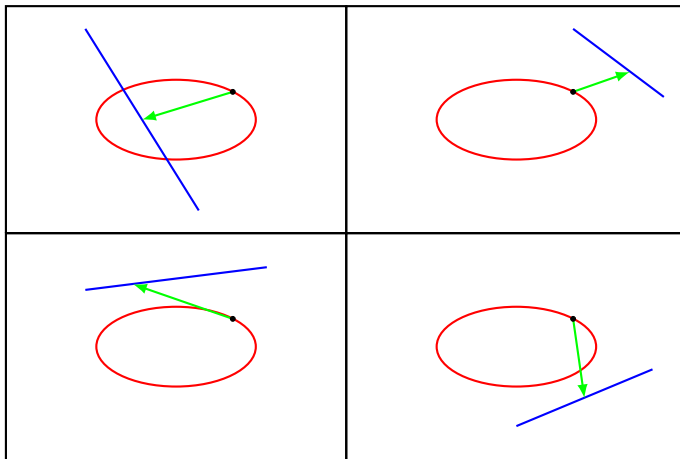


Figure: Illustrations of SQP subproblem solutions

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

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

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

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

and it is easily shown that

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

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: **compute step**: solve

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} d_k \\ y_k \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) \\ c_k \end{bmatrix}$$

4: **update merit parameter**: set τ_k to ensure $\Delta q(x_k, \tau_k, \nabla f(x_k), d_k) \gg 0$, offered by

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

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

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

6: **end for**

Convergence theory

Assumption

- ▶ $f, c, \nabla f$, 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 $k \in \mathbb{N}$

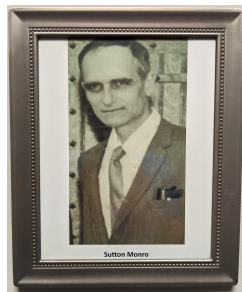
Theorem SQP-B

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

$$\|d_k\|_2 \rightarrow 0, \quad \|c_k\|_2 \rightarrow 0, \quad \|\nabla f(x_k) + J_k^T y_k\|_2 \rightarrow 0$$

Stochastic gradient method (SG)

Invented by Herbert Robbins and Sutton Monro (1951)



Sutton Monro, former Lehigh faculty member

Stochastic gradient (*not* descent)

Consider the stochastic optimization problem

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

where $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ 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 $x_{k+1} \leftarrow x_k - \alpha_k g_k$, where $\mathbb{E}_k[g_k] = \nabla f(x_k)$ and $\mathbb{E}_k[\|g_k - \nabla f(x_k)\|_2^2] \leq M$
 - 4: **end for**
-

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

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

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

SG illustration

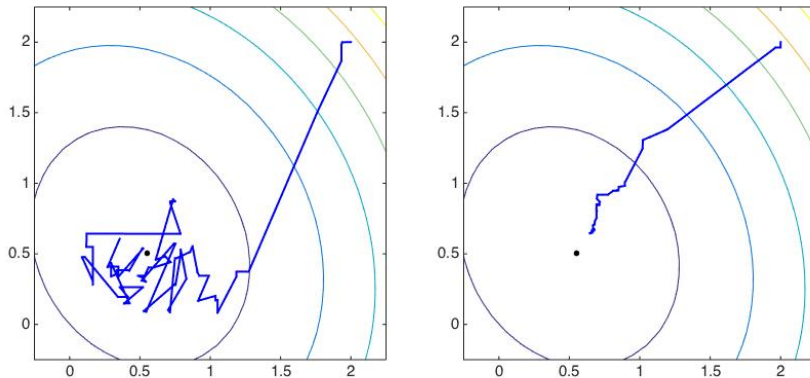


Figure: SG with fixed step size (left) vs. diminishing step sizes (right)

SG theory

Theorem SG

Since $\mathbb{E}_k[g_k] = \nabla f(x_k)$ and $\mathbb{E}_k[\|g_k - \nabla f(x_k)\|_2^2] \leq M$ for all $k \in \mathbb{N}$:

$$\begin{aligned} \alpha_k = \frac{1}{L} &\implies \mathbb{E} \left[\frac{1}{k} \sum_{j=1}^k \|\nabla f(x_j)\|_2^2 \right] \leq \mathcal{O}(M) \\ \alpha_k = \Theta \left(\frac{1}{k} \right) &\implies \mathbb{E} \left[\frac{1}{\left(\sum_{j=1}^k \alpha_j \right)} \sum_{j=1}^k \alpha_j \|\nabla f(x_j)\|_2^2 \right] \rightarrow 0 \\ &\implies \liminf_{k \rightarrow \infty} \mathbb{E}[\|\nabla f(x_k)\|_2^2] = 0 \end{aligned}$$

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

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

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

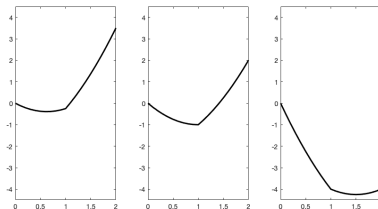


Figure: Three cases for upper bound of ϕ

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: **compute step**: solve

$$\begin{bmatrix} H_k & J_k^T \\ J_k & 0 \end{bmatrix} \begin{bmatrix} d_k \\ y_k \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) \\ c_k \end{bmatrix}$$

4: **update merit parameter**: set τ_k to ensure $\Delta q(x_k, \tau_k, \nabla f(x_k), d_k) \gg 0$, offered by

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

5: **compute step size**: set

$$\hat{\alpha}_k \leftarrow \frac{2(1 - \eta) \Delta q(x_k, \tau_k, \nabla f(x_k), d_k)}{(\tau_k L_k + \Gamma_k) \|d_k\|_2^2} \quad \text{and} \quad \tilde{\alpha}_k \leftarrow \hat{\alpha}_k - \frac{4 \|c_k\|_1}{(\tau_k L_k + \Gamma_k) \|d_k\|_2^2}$$

6: **then**

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

7: **then** 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.*

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 g_k with

$$\mathbb{E}_k[g_k] = \nabla f(x_k) \quad \text{and} \quad \mathbb{E}_k[\|g_k - \nabla f(x_k)\|_2^2] \leq M$$

Search directions computed by:

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

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

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

3: **compute step**: 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: **update merit parameter**: set τ_k to ensure $\Delta q(x_k, \tau_k, g_k, d_k) \gg 0$, 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: **compute step size**: set

$$\hat{\alpha}_k \leftarrow \frac{\beta_k \Delta q(x_k, \tau_k, g_k, d_k)}{(\tau_k L + \Gamma)\|d_k\|_2^2} \quad \text{and} \quad \tilde{\alpha}_k \leftarrow \hat{\alpha}_k - \frac{4\|c_k\|_1}{(\tau_k L + \Gamma)\|d_k\|_2^2}$$

6: **then**

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

7: **then** $x_{k+1} \leftarrow x_k + \alpha_k d_k$

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 $\hat{\alpha}_k$ and $\tilde{\alpha}_k$ onto

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

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

Fundamental lemmas

Lemma

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

$$\begin{aligned} & \phi(x_k + \alpha_k d_k, \tau_k) - \phi(x_k, \tau_k) \\ & \leq \underbrace{-\alpha_k \Delta q(x_k, \tau_k, \nabla f(x_k), d_k^{\text{true}})}_{\mathcal{O}(\beta_k), \text{"deterministic"}} + \underbrace{\frac{1}{2} \alpha_k \beta_k \Delta q(x_k, \tau_k, g_k, d_k)}_{\mathcal{O}(\beta_k^2), \text{stochastic/noise}} + \underbrace{\alpha_k \tau_k \nabla f(x_k)^T (d_k - d_k^{\text{true}})}_{\text{due to adaptive } \alpha_k} \end{aligned}$$

Lemma

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

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

as well as

$$\begin{aligned} \nabla f(x_k)^T d_k^{\text{true}} & \geq \mathbb{E}_k[g_k^T d_k] \geq (\nabla f(x_k)^T d_k)^{\text{true}} - \zeta^{-1} M \quad \text{and} \\ \mathbb{E}_k[d_k^T H_k d_k] & \geq d_k^{\text{true}T} H_k d_k^{\text{true}} \end{aligned}$$

Good merit parameter behavior

Lemma

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

$$\mathbb{E}_k[\alpha_k \tau_k \nabla f(x_k)^T (d_k - d_k^{\text{true}})] = \beta_k^2 \tau_{\min} \mathcal{O}(\sqrt{M})$$

Theorem

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

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

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

Good merit parameter behavior

Lemma

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

$$\mathbb{E}_k[\alpha_k \tau_k \nabla f(x_k)^T (d_k - d_k^{\text{true}})] = \beta_k^2 \tau_{\min} \mathcal{O}(\sqrt{M})$$

Theorem

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

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

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

Poor merit parameter behavior

$\{\tau_k\} \searrow 0$:

- ▶ cannot occur if $\|g_k - \nabla f(x_k)\|_2$ is bounded uniformly
- ▶ occurs with small probability if distribution of g_k has *fast decay*

$\{\tau_k\}$ remains too large:

- ▶ if there exists $p \in (0, 1]$ such that, for all k in infinite \mathcal{K} ,

$$\mathbb{P}_k \left[g_k^T d_k + \max\{d_k^T H_k d_k, 0\} \geq \nabla f(x_k)^T d_k^{\text{true}} + \max\{(d_k^{\text{true}})^T H_k d_k^{\text{true}}, 0\} \right] \geq p$$

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

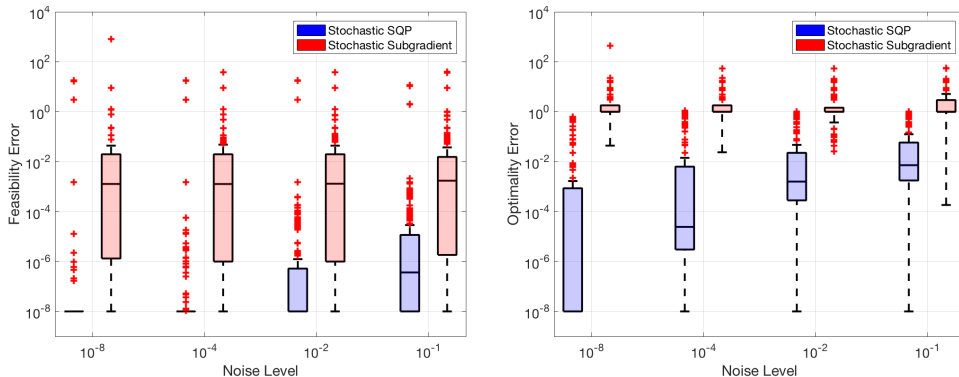


Figure: Box plots for feasibility errors (left) and optimality errors (right).

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

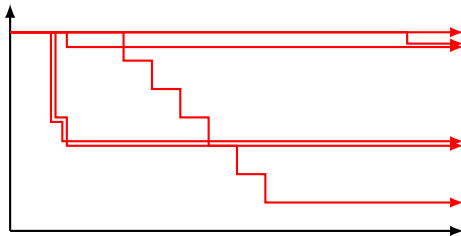
$$\left(\frac{\tau_{-1}(f_0 - f_{\inf}) + \|c_0\|_1}{\min\{\kappa_1, \kappa_2 \tau_{\min}\}} \right) \epsilon^{-2}.$$

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 .



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

Theorem 9

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 10

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

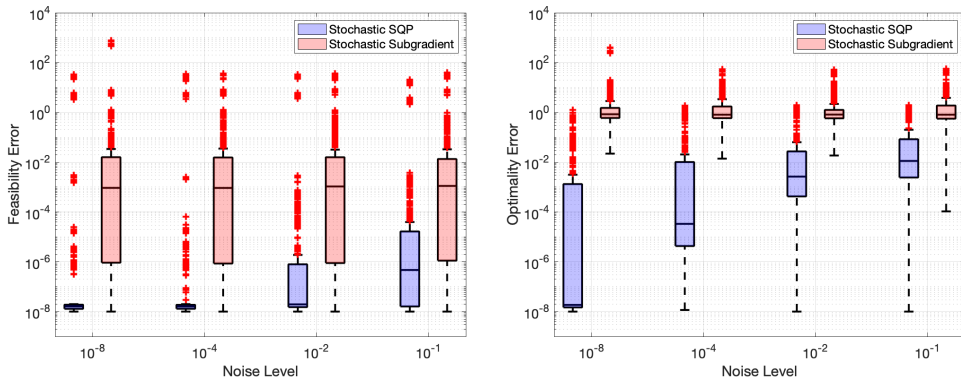


Figure: Box plots for feasibility errors (left) and optimality errors (right).

Recent work (under review): Matrix-free methods

Inexact subproblem solves

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

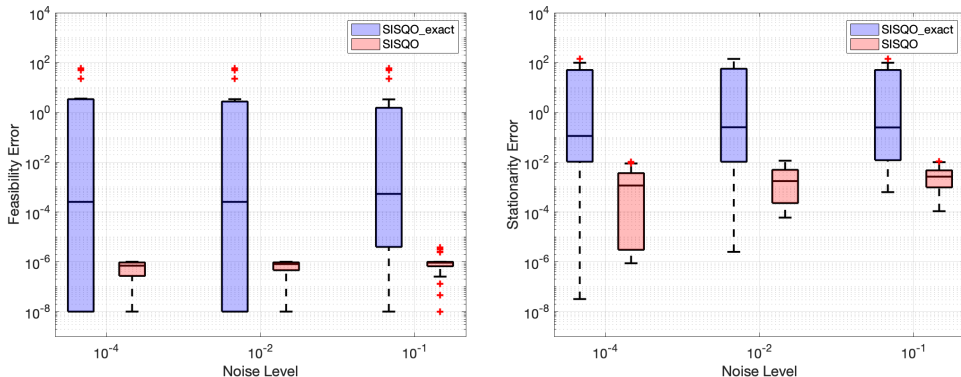


Figure: Box plots for feasibility errors (left) and optimality errors (right).

Current work: Inequality constraints

Inequality constraints

- ▶ SQP
- ▶ interior-point

Main challenge: For *equalities* 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

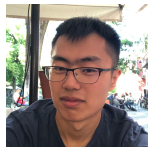
Summary

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>.