# Stochastic-Gradient-based Interior-Point Methods

**Frank E. Curtis**, Lehigh University

presented at

60th Annual Allerton Conference on Communication, Control, and Computing

September 26, 2024

# Collaborators and references



Submitted papers:

▶ F. E. Curtis, V. Kungurtsev, D. P. Robinson, and Q. Wang, "A Stochastic-Gradient-based Interior-Point Algorithm for Solving Smooth Bound-Constrained Optimization Problems," https://arxiv.org/abs/2304.14907, in third round of review (SIAM Journal on Optimization).

▶ F. E. Curtis, X. Jiang, and Q. Wang, "Single-Loop Deterministic and Stochastic Interior-Point Algorithms for Nonlinearly Constrained Optimization," https://arxiv.org/abs/2408.16186, in first round of review (Mathematical Programming, Series B).

# Outline

# Outline

# Motivation

Interior-point methods are the workhorse for deterministic nonlinearly constrained optimization.

- ▶ Ipopt, Knitro, LOQO, etc.

Before our work, there were no stochastic interior-point methods with convergence guarantees.[†]

Why not?

- ▶ Stochastic algorithms for constrained optimization are not widely studied
- ▶ ...except for projection methods, manifold-based methods, and conditional gradient methods.
- ▶ Stochastic-gradient-based algorithms require gradients to be bounded and Lipschitz continuous
- ▶ ...but barrier functions (e.g., logarithmic barrier) have neither property.

In our first paper and this talk, we focus on the bound-constrained case.

- ▶ I will end with the additional discussion about the generally constrained case.

---

[†]An idea was proposed, but there was a flaw in the analysis.

## Bound-constrained setting

Given $f : \mathbb{R}^n \to \mathbb{R}$ and $(l, u) \in \mathbb{R}^n \times \mathbb{R}^n$ with $l < u$, consider

$$\min_{x \in \mathbb{R}^n} \ f(x)$$
$$\text{s.t. } l \leq x \leq u$$

If $x$ is a minimizer, then for some $(y, z)$ one has

$$\nabla f(x) - y + z = 0, \ \ 0 \leq (x - l) \perp y \geq 0, \ \ 0 \leq (u - x) \perp z \geq 0.$$

(We can handle infinite bounds, but in this talk consider finite bounds for simplicity....)

## Textbook algorithm

For all $\mu \in \mathbb{R}_{>0}$, consider the barrier-augmented function

$$\phi(x, \mu) = f(x) - \mu \sum_{i=1}^{n} \log(x_i - l_i) - \mu \sum_{i=1}^{n} \log(u_i - x_i).$$

---

**Algorithm IPM** : Interior-point method (textbook version)

---

1: choose an initial point $x_1 \in (l, u)$ and barrier parameter $\mu_0 \in \mathbb{R}_{>0}$
2: **for** all $k \in \{1, 2, \dots\}$ **do**
3:     **if** $\|\nabla_x \phi(x_k, \mu_{k-1})\|_2 \leq \theta \mu_{k-1}$ **then** set $\mu_k \leq \mu_{k-1}$ **else** set $\mu_k \leftarrow \mu_{k-1}$
4:     compute descent direction $d_k$ (e.g., $-\nabla \phi(x_k, \mu_k)$)
5:     set $\alpha_{k,\max} \in (0, 1]$ by fraction-to-the-boundary rule to ensure

$$x_k + \alpha_{k,\max} d_k - l \geq \epsilon(x_k - l) \quad \text{and} \quad u - (x_k + \alpha_{k,\max} d_k) \geq \epsilon(u - x_k)$$

6:     set $\alpha_k \in (0, \alpha_{k,\max}]$ to ensure sufficient decrease $\phi(x_{k+1}, \mu_k) \ll \phi(x_k, \mu_k)$
7: **end for**

---

**Note**: Essentially a nested-loop algorithm with inner loop having fixed $\mu$
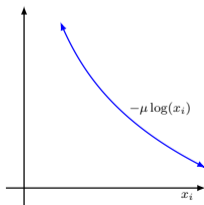
## Major challenges for the stochastic setting

Stationarity test:

- ▶ Computing $\|\nabla_x \phi(x_k, \mu_{k-1})\|_2$ is intractable
- ▶ Could estimate it using a stochastic gradient, but then a probabilistic guarantee, at best

Fraction-to-the-boundary rule:

- ▶ Tying fraction to current iterate $x_k$ leads to issues
- ▶ ... stochastic gradients could push iterate sequence to boundary too quickly

Unbounded gradients and lack of Lipschitz continuity:

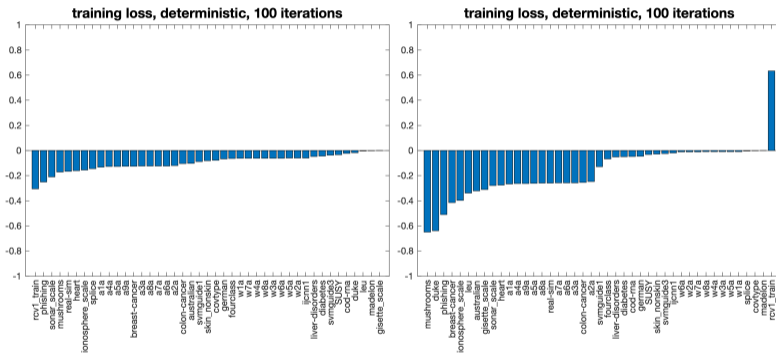## Our approach

Our approach is based on two coupled ideas:

▶ prescribed decreasing barrier parameter sequence $\{\mu_k\} \searrow 0$ (single-loop algorithm!)

▶ prescribed $\{\theta_k\} \searrow 0$ and enforcement of

$$x_{k+1} \in \mathcal{N}_{[l,u]}(\theta_k) := \{x \in \mathbb{R}^n : l + \theta_k \leq x \leq u - \theta_k\}$$

"Wait! Is it worthwhile to have an algorithm like this?!"

▶ Our experiments say yes!

## Deterministic setting



Relative performance of SLIP vs. PGM, deterministic setting, training logistic regression (left) and neural network models with one hidden layer with cross-entropy loss (right).

## Proposed algorithm

---

**Algorithm SLIP** : Single-loop interior-point method

---

1: choose an initial point $x_1 \in \mathcal{N}_{[l,u]}(\theta_0)$, $\{\mu_k\} \searrow 0$, $\{\theta_k\} \searrow 0$
2: **for** all $k \in \{1, 2, \dots\}$ **do**
3:     compute descent direction $d_k$ (e.g., estimating $-\nabla\phi(x_k, \mu_k)$)
4:     set

$$\alpha_k \leftarrow \frac{1}{L + 2\mu_k \theta_k^{-2}}$$

5:     set $\gamma_k \in (0, 1]$ to ensure

$$x_{k+1} \leftarrow x_k + \gamma_k \alpha_k d_k \in \mathcal{N}_{[l,u]}(\theta_k)$$

6: **end for**

---

**Note**: Our paper considers a more general framework; this is a simplified instance

# Key observation

Our first key observation is that the algorithm essentially acts equivalently to minimize

$$\phi(x, \mu) = f(x) - \mu \sum_{i=1}^{n} \log(x_i - l_i) - \mu \sum_{i=1}^{n} \log(u_i - x_i)$$

and

$$\tilde{\phi}(x, \mu) = f(x) - \mu \sum_{i=1}^{n} \log\left(\frac{x_i - l_i}{\chi}\right) - \mu \sum_{i=1}^{n} \log\left(\frac{u_i - x_i}{\chi}\right),$$

where $\chi$ is sufficiently large such that $\frac{x_i - l_i}{\chi} \in [0, 1]$ and $\frac{u_i - x_i}{\chi} \in [0, 1]$ for all $i \in [n]$.

The latter is simply a shifted form of the other.

▶ They have the same gradients! $\nabla_x \phi(x, \mu) = \nabla_x \tilde{\phi}(x, \mu)$
▶ For the latter, $\bar{\mu} < \mu$ implies that $\tilde{\phi}(x, \bar{\mu}) < \tilde{\phi}(x, \mu)$.

The algorithm uses $\phi$, but our analysis can focus on monotonically decreasing $\{\tilde{\phi}(x_k, \mu_k)\}$.

## Critical lemmas, deterministic setting

**Lemma**

*For all $k \in \mathbb{N}$, one finds for $L_k := L + 2\mu_k \theta_k^{-2}$ that*

$$\tilde{\phi}(x_{k+1}, \mu_k) \leq \tilde{\phi}(x_k, \mu_k) + \nabla_x \tilde{\phi}(x_k, \mu_k)^T (x_{k+1} - x_k) + \tfrac{1}{2} L_k \|x_{k+1} - x_k\|_2^2,$$

*so $\{\alpha_k\} = \{L_k^{-1}\} \implies \tilde{\phi}(x_{k+1}, \mu_{k+1}) \leq \tilde{\phi}(x_k, \mu_k) - \tfrac{1}{2}\gamma_k \alpha_k \|\nabla_x \tilde{\phi}(x_k, \mu_k)\|_2^2.$*

**Lemma**

*For all $k \in \mathbb{N}$, one finds that $\gamma_k$ is bounded below by the minimum of 1 and*

$$\alpha_k^{-1} \left( \frac{\tfrac{1}{2}\mu_k \Delta}{\mu_k + \tfrac{1}{2}\kappa_{\nabla f} \Delta} - \theta_k \right) (\kappa_{\nabla f} + \mu_k \theta_{k-1}^{-1})^{-1}.$$

*Thus, with $t \in [-1, 0)$, $\{\mu_k\} = \{\mu_1 k^t\}$, $\{\theta_{k-1}\} = \{\theta_0 k^t\}$, and $\{\alpha_k\} = \{L_k^{-1}\}$, one finds that*

$$\sum_{k=1}^{\infty} \gamma_k \alpha_k = \infty \quad \text{and} \quad \{\mu_k \theta_{k-1}^{-1}\} \text{ is bounded.}$$

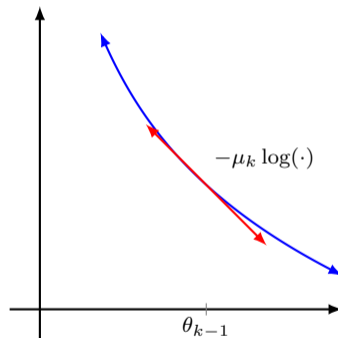# Convergence guarantee, deterministic setting

## Theorem

*One finds that*
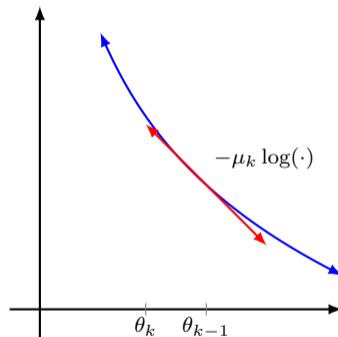
$$\liminf_{k \to \infty} \|\nabla_x \phi(x_k, \mu_k)\|_2^2 = 0,$$

*and, for any infinite-cardinality set $\mathcal{K} \subseteq \mathbb{N}$ such that $\{\nabla_x \phi(x_k, \mu_k)\}_{k \in \mathcal{K}} \to 0$ and $\{x_k\}_{k \in \mathcal{K}} \to \bar{x}$, the limit point $\bar{x}$ is a KKT point (i.e., there exists $\bar{y}$ and $\bar{z}$ such that $(\bar{x}, \bar{y}, \bar{z})$ satisfies KKT conditions).*
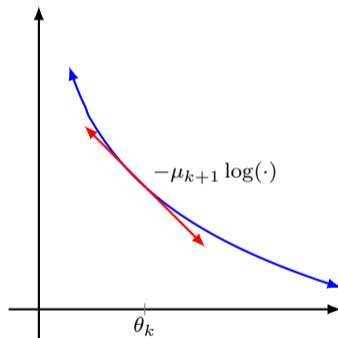
## Why does it work?

# Why does it work?

## Why does it work?

# Outline

Single-Loop Interior-Point (SLIP) Method    **Stochastic Bound-Constrained Setting**    Generally Constrained Setting    Conclusion

○○○○○○○○○○○○     ○●○○○○○○○○○      ○○○○      ○○○

## Stochastic setting

In the stochastic setting, the algorithm parameters need to be chosen more carefully!

▶ Notably, $\gamma_k$ needs to be chosen based on knowledge of noise bound.

▶ For the deterministic setting, $\{\mu_k\} = \{\mu_1 k^t\}$ and $\{\theta_{k-1}\} = \{\theta_0 k^t\}$ for $t = -1$ implies

$$\{\alpha_k\} = \left\{ \frac{1}{L + 2\mu_k \theta_k^{-2}} \right\} = \Theta(k^t),$$

but for stochastic setting, step-size sequence $\{\alpha_k\}$ can no longer decrease at same rate as $\{\mu_k\}$.

▶ It needs to decrease more slowly than $\{\mu_k\}$ (although rates can be arbitrarily close).

## Accounting for the error

The issue arises from the following lemma.

> **Lemma**
>
> *For all $k \in \mathbb{N}$, one finds that*
>
> $$\tilde{\phi}(X_{k+1}, \mu_{k+1}) - \tilde{\phi}(X_k, \mu_k)$$
> $$\leq -\Gamma_k A_k \|\nabla_x \tilde{\phi}(X_k, \mu_k)\|_{H_k^{-1}}^2 + \Gamma_k A_k \nabla_x \tilde{\phi}(X_k, \mu_k)^T H_k^{-1}(\nabla_x \tilde{\phi}(X_k, \mu_k) - Q_k)$$
> $$+ \tfrac{1}{2}\Gamma_k^2 A_k^2 \lambda_{k,\min}^{-1} \ell_{\nabla f, \mathcal{B}, k} \|Q_k\|_{H_k^{-1}}^2.$$

Using $\{\mu_k\} = \{\mu_1 k^{-1}\}$ and $\{\theta_{k-1}\} = \{\theta_0 k^{-1}\}$, so $\{\alpha_k\} = \Theta(k^t)$, leaves the final term uncontrolled!
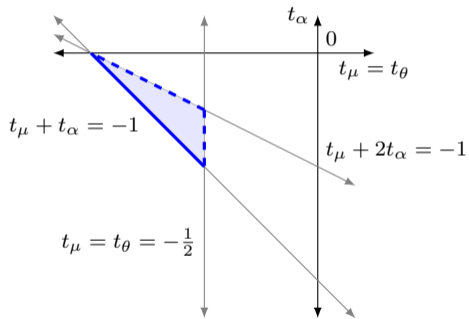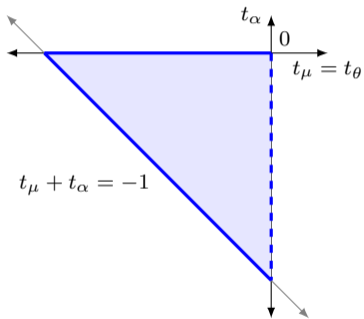
## Parameter rule

Given prescribed $(t_\mu, t_\theta, t_\alpha) \in (-\infty, -\frac{1}{2}) \times (-\infty, -\frac{1}{2}) \times (-\infty, 0)$ such that $t_\mu = t_\theta$, $t_\mu + t_\alpha \in [-1, 0)$, and $t_\mu + 2t_\alpha \in (-\infty, -1)$ along with prescribed $\alpha_{\text{buff}} \in \mathbb{R}_{\geq 0}$, $\{\alpha_{k,\text{buff}}\} \subset \mathbb{R}_{\geq 0}$, $\gamma_{\text{buff}} \in \mathbb{R}_{\geq 0}$, and $\{\gamma_{k,\text{buff}}\} \subset \mathbb{R}_{\geq 0}$ such that $\alpha_{k,\text{buff}} \leq \alpha_{\text{buff}} k^{2t_\mu}$ and $\gamma_{k,\text{buff}} \leq \gamma_{\text{buff}} k^{t_\mu}$ for all $k \in \mathbb{N}$, the algorithm employs

$$\alpha_{k,\min} := \frac{\lambda_{k,\min} k^{t_\alpha}}{\ell_{\nabla f, \mathcal{B}} + 2\mu_k \theta_k^{-2}}, \qquad \gamma_{k,\min} := \min\left\{1, \frac{\lambda_{k,\min}\left(\frac{\frac{1}{2}\mu_k \Delta}{\mu_k + \frac{1}{2}(\kappa_{\nabla f, \mathcal{B}, \infty} + \sigma_\infty)\Delta} - \theta_k\right)}{\alpha_{k,\max}(\kappa_{\nabla f, \mathcal{B}, \infty} + \sigma_\infty + \mu_k \theta_{k-1}^{-1})}\right\},$$

$$\alpha_{k,\max} := \alpha_{k,\min} + \alpha_{k,\text{buff}}, \qquad \text{and} \quad \gamma_{k,\max} := \min\{1, \gamma_{k,\min} + \gamma_{k,\text{buff}}\}$$

and makes a (run-and-iterate-dependent) choice $\alpha_k \in \min\left\{\frac{\lambda_{k,\min} k^{t_\alpha}}{L + 2\mu_k \theta_k^{-2}}, \alpha_{k,\max}\right\}$ for all $k \in \mathbb{N}$.

## Acceptable rate values



$t_\alpha$
$0$
$t_\mu = t_\theta$

$t_\mu + t_\alpha = -1$

$t_\alpha$
$0$
$t_\mu = t_\theta$

$t_\mu + t_\alpha = -1$

$t_\mu + 2t_\alpha = -1$

$t_\mu = t_\theta = -\frac{1}{2}$

## Convergence guarantee, stochastic setting

**Theorem**

Suppose $t \in (-1, -\frac{1}{2})$ and $t_\alpha \in (-\infty, 0)$ have

$$t + t_\alpha \in [-1, 0) \quad and \quad t + 2t_\alpha \in (-\infty, -1)$$

and for some $\sigma \in \mathbb{R}_{>0}$ one has for all $k \in \mathbb{N}$ that

$$\mathbb{E}[G_k | \mathcal{F}_k] = \nabla f(X_k) \quad and \quad \|G_k - \nabla f(X_k)\|_2 \leq \sigma.$$

Then, with $\{\mu_k\} = \{\mu_1 k^t\}$, $\{\theta_{k-1}\} = \{\theta_0 k^t\}$, and $\{\alpha_k\} = \{L_k^{-1} k^{t_\alpha}\}$, one finds that

$$\liminf_{k \to \infty} \|\nabla_x \phi(X_k, \mu_k)\|_2^2 = 0 \quad almost\ surely.$$

Consequently, considering any realization $\{x_k\}$ of $\{X_k\}$, for any infinite-cardinality set $\mathcal{K} \subseteq \mathbb{N}$ such that $\{\nabla_x \phi(x_k, \mu_k)\}_{k \in \mathcal{K}} \to 0$ and $\{x_k\}_{k \in \mathcal{K}} \to \bar{x}$, the limit point $\bar{x}$ is a KKT point.

Single-Loop Interior-Point (SLIP) Method    **Stochastic Bound-Constrained Setting**    Generally Constrained Setting    Conclusion

○○○○○○○○○○○○    ○○○○○○●○○○    ○○○○    ○○○

# Numerical experiments

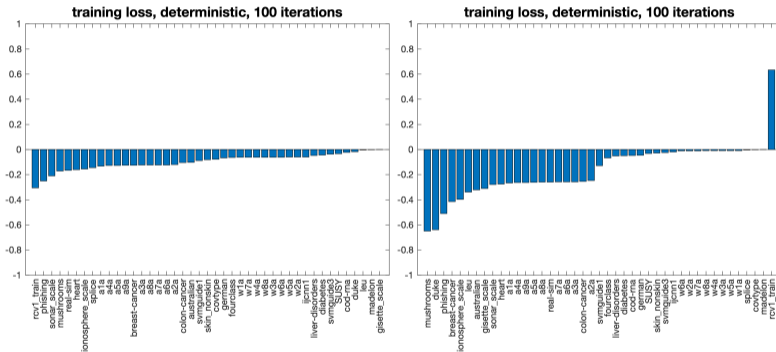Compare SLIP with a projected stochastic gradient method (PSGM) for which

$$x_{k+1} \leftarrow \text{Proj}_{[l,u]}(x_k + \alpha_k d_k).$$

Experiments involve:

- ▶ binary classification problems with LIBSVM datasets
- ▶ two classifiers:
  - ▶ logistic regression (convex) and
  - ▶ neural network with one hidden layer and cross-entropy loss (nonconvex)
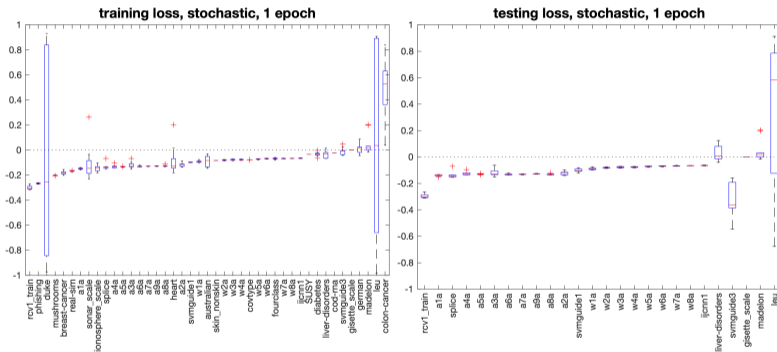- ▶ performance measure

$$\frac{f(x_{\text{end}}^{\text{SLIP}}) - f(x_{\text{end}}^{\text{PSGM}})}{\max\{f(x_{\text{end}}^{\text{SLIP}}), f(x_{\text{end}}^{\text{PSGM}}), 1\}} \in (-1, 1)$$
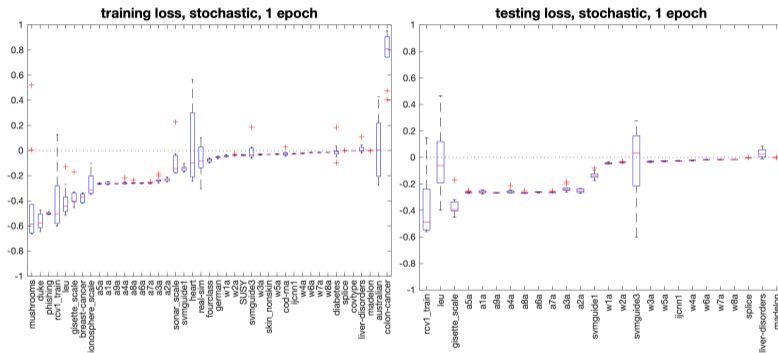
## Deterministic setting



Relative performance of SLIP and PGM, deterministic setting, training logistic regression (left) and neural network models with one hidden layer with cross-entropy loss (right).

## Stochastic setting, logistic regression



Relative performance of SLIP and PSGM, stochastic setting (10 runs each), training logistic regression models; among 43 training datasets, 26 have testing datasets.

Single-Loop Interior-Point (SLIP) Method    **Stochastic Bound-Constrained Setting**    Generally Constrained Setting    Conclusion

○○○○○○○○○○○○    ○○○○○○○○○○●    ○○○○    ○○○

## Stochastic setting, neural network with cross-entropy loss



Relative performance of SLIP and PSGM, stochastic setting (10 runs each), training neural network models (with one hidden layer) with cross-entropy loss; among 43 training datasets, 26 have testing datasets.

# Outline

## SLIP algorithm

---

**Algorithm SLIP** : Single-loop interior-point method

---

1: choose an initial point $x_1 \in \mathcal{N}_{[l,u]}(\theta_0)$, $\{\mu_k\} \searrow 0$, $\{\theta_k\} \searrow 0$
2: **for** all $k \in \{1, 2, \dots\}$ **do**
3:      compute descent direction $d_k$ (e.g., estimating $-\nabla\phi(x_k, \mu_k)$)
4:      set

$$\alpha_k \leftarrow \frac{1}{L + 2\mu_k \theta_k^{-2}}$$

5:      set $\gamma_k \in (0, 1]$ to ensure

$$x_{k+1} \leftarrow x_k + \gamma_k \alpha_k d_k \in \mathcal{N}_{[l,u]}(\theta_k)$$

6: **end for**

---

How can this be extended for the generally constrained setting?

- ▶ This is a feasible algorithm.
- ▶ Neighborhood enforcement is the real issue! Constraint value depends nonlinearly on $\gamma_k$.

## Search direction conditions

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

$$\phi(x, \mu) = f(x) - \mu \sum_{i=1}^n \log(-c_i(x))$$
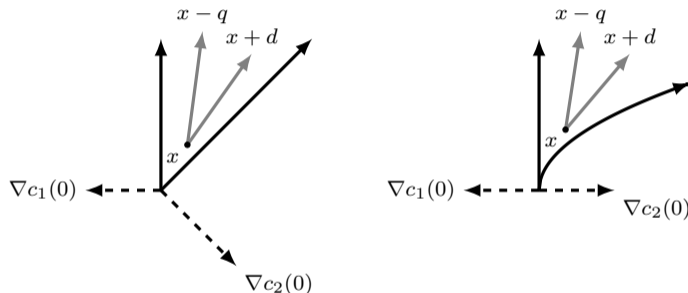
Need an initial point $x_1 \in \mathbb{R}^n$ satisfying

$$Ax_1 = b \ \text{ and } \ c(x_1) < 0,$$

and, with $P := I - A^T(AA^T)^{-1}A$, to ensure/assume that, for all $k \in \mathbb{N}$, one can compute $d_k$ satisfying

$$Ad_k = 0$$
$$\underline{\zeta}\|Pq_k\|_2 \leq \|d_k\|_2 \leq \overline{\zeta}\|Pq_k\|_2$$
$$-(Pq_k)^T d_k \geq \zeta\|Pq_k\|_2\|d_k\|_2$$
$$\nabla c_i(x_k)^T d_k \leq -\tfrac{1}{2}\overline{\eta}\|d_k\|_2 \ \text{ for all } \ i \in \{j \in [m] : -\eta\mu_k < c_i(x_k)\}.$$

# Main challenge



Assuming nice conditions (e.g., on the left, not on the right) and parameter choices similar to the bound-constrained case, we prove that the projected gradient of the barrier-augmented function vanishes and, if a limit point satisfies the LICQ, then the limit point is a KKT point.

# Outline

## Summary

Presented a single-loop interior-point method for solving bound-constrained problems, with

- ▶ prescribed barrier and "neighborhood" parameter sequences,
- ▶ no need for stationarity tests, fraction-to-the-boundary rules, or line searches,
- ▶ convergence guarantees in deterministic and stochastic settings, and
- ▶ promising numerical performance!

Presented an overview of our extension to the "generally constrained" setting.

- ▶ There is more to be done!

## Collaborators and references



Submitted papers:

▶ F. E. Curtis, V. Kungurtsev, D. P. Robinson, and Q. Wang, "A Stochastic-Gradient-based Interior-Point Algorithm for Solving Smooth Bound-Constrained Optimization Problems," https://arxiv.org/abs/2304.14907, in third round of review (SIAM Journal on Optimization).

▶ F. E. Curtis, X. Jiang, and Q. Wang, "Single-Loop Deterministic and Stochastic Interior-Point Algorithms for Nonlinearly Constrained Optimization," https://arxiv.org/abs/2408.16186, in first round of review (Mathematical Programming, Series B).