Krylov Cubic Regularized Newton: A Subspace Second-Order Method with Dimension-Free Convergence Rate

Ruichen Jiang UT Austin Parameswaran Raman AWS

Shoham Sabach Technion & AWS Aryan Mokhtari UT Austin

Mingyi Hong UMN & AWS Volkan Cevher EPFL & AWS

January 12, 2024



1 Introduction

The Proposed Method

③ Convergence Analysis

4 Numerical Results

Consider the unconstrained minimization problem

 $\min_{\boldsymbol{x}\in\mathbb{R}^d}f(\boldsymbol{x}),$

where f is convex and twice continuously differentiable

Consider the unconstrained minimization problem

 $\min_{\boldsymbol{x}\in\mathbb{R}^d}f(\boldsymbol{x}),$

where f is convex and twice continuously differentiable

Popular methods: first-order methods such as gradient descent

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$$

- Cheap to implement
- Slow convergence, esp. for ill-conditioned problems X

- ▶ We focus on second-order methods that utilize the Hessian of *f*
- Newton's method

$$oldsymbol{x}_{k+1} = oldsymbol{x}_k -
abla^2 f(oldsymbol{x}_k)^{-1}
abla f(oldsymbol{x}_k)$$

- We focus on second-order methods that utilize the Hessian of f
- Newton's method

$$oldsymbol{x}_{k+1} = oldsymbol{x}_k -
abla^2 f(oldsymbol{x}_k)^{-1}
abla f(oldsymbol{x}_k)$$

Cubic regularized Newton (CRN) method [Griewank'81; Nesterov-Polyak'06]

$$\boldsymbol{x}_{k+1} = \operatorname*{argmin}_{\boldsymbol{x} \in \mathbb{R}^d} \left\{ f(\boldsymbol{x}_k) + \nabla f(\boldsymbol{x}_k)^\top (\boldsymbol{x} - \boldsymbol{x}_k) + \frac{1}{2} (\boldsymbol{x} - \boldsymbol{x}_k)^\top \nabla^2 f(\boldsymbol{x}_k) (\boldsymbol{x} - \boldsymbol{x}_k) + \frac{M}{6} \| \boldsymbol{x} - \boldsymbol{x}_k \|^3 \right\}$$

- We focus on second-order methods that utilize the Hessian of f
- Newton's method

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k) \\ &= \operatorname*{argmin}_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^\top \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k) \right\} \end{aligned}$$

Cubic regularized Newton (CRN) method [Griewank'81; Nesterov-Polyak'06]

$$oldsymbol{x}_{k+1} = rgmin_{oldsymbol{x}\in\mathbb{R}^d} \left\{ f(oldsymbol{x}_k) +
abla f(oldsymbol{x}_k)^ op (oldsymbol{x} - oldsymbol{x}_k) + rac{1}{2} (oldsymbol{x} - oldsymbol{x}_k)^ op
abla^2 f(oldsymbol{x}_k) (oldsymbol{x} - oldsymbol{x}_k) + rac{M}{6} \|oldsymbol{x} - oldsymbol{x}_k\|^3
ight\}$$

- We focus on second-order methods that utilize the Hessian of f
- Newton's method

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k) \\ &= \operatorname*{argmin}_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^\top \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k) \right\} \end{aligned}$$

Cubic regularized Newton (CRN) method [Griewank'81; Nesterov-Polyak'06]

$$oldsymbol{x}_{k+1} = rgmin_{oldsymbol{x}\in\mathbb{R}^d} \left\{ f(oldsymbol{x}_k) +
abla f(oldsymbol{x}_k)^ op (oldsymbol{x} - oldsymbol{x}_k) + rac{1}{2} (oldsymbol{x} - oldsymbol{x}_k)^ op
abla^2 f(oldsymbol{x}_k) (oldsymbol{x} - oldsymbol{x}_k) + rac{M}{6} \|oldsymbol{x} - oldsymbol{x}_k\|^3
ight\}$$

• When f is convex,
$$f(\mathbf{x}_k) - f^* = \mathcal{O}(1/k^2)$$

• When f is strongly convex, it achieves a superlinear convergence rate

▶ The main drawback of CRN is its substantial memory and computational costs

- Computing & storing the Hessian $\nabla^2 f(\mathbf{x})$: $\mathcal{O}(d^2)$
- Solving a cubic subproblem: $\mathcal{O}(d^3)$

▶ As a result, CRN becomes impractical for optimization problems with high dimensions

► To reduce the cost, one approach is to execute 2nd-order updates in a subspace V_k of dimension m ≪ d [Doikov-Richtárik'18; Gower et al.'19; Hanzely et al.'20]

► To reduce the cost, one approach is to execute 2nd-order updates in a subspace V_k of dimension m ≪ d [Doikov-Richtárik'18; Gower et al.'19; Hanzely et al.'20]

• Let
$$\boldsymbol{g}_k :=
abla f(\boldsymbol{x}_k)$$
 and $\boldsymbol{\mathsf{H}}_k :=
abla^2 f(\boldsymbol{x}_k)$

CRN:
$$egin{array}{ccc} oldsymbol{s}_k = \operatorname*{argmin}_{oldsymbol{s} \in \mathbb{R}^d} \left\{ oldsymbol{g}_k^{ op} oldsymbol{s} + rac{1}{2} oldsymbol{s}^{ op} oldsymbol{H}_k oldsymbol{s} + rac{M}{6} \|oldsymbol{s}\|^3
ight\} \ \end{array}$$

$$oldsymbol{x}_{k+1} = oldsymbol{x}_k + oldsymbol{s}_k$$

(

► To reduce the cost, one approach is to execute 2nd-order updates in a subspace V_k of dimension m ≪ d [Doikov-Richtárik'18; Gower et al.'19; Hanzely et al.'20]

. .

• Let
$$\boldsymbol{g}_k := \nabla f(\boldsymbol{x}_k)$$
 and $\boldsymbol{\mathsf{H}}_k := \nabla^2 f(\boldsymbol{x}_k)$

Subspace CRN:
$$\boldsymbol{s}_{k} = \underset{\boldsymbol{s} \in \boldsymbol{\mathcal{V}}_{k}}{\operatorname{argmin}} \left\{ \boldsymbol{g}_{k}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \boldsymbol{\mathsf{H}}_{k} \boldsymbol{s} + \frac{M}{6} \| \boldsymbol{s} \|^{3} \right\}$$

$$oldsymbol{x}_{k+1} = oldsymbol{x}_k + oldsymbol{s}_k$$

► To reduce the cost, one approach is to execute 2nd-order updates in a subspace V_k of dimension m ≪ d [Doikov-Richtárik'18; Gower et al.'19; Hanzely et al.'20]

• Let
$$\boldsymbol{g}_k :=
abla f(\boldsymbol{x}_k)$$
 and $\boldsymbol{\mathsf{H}}_k :=
abla^2 f(\boldsymbol{x}_k)$

Subspace CRN:

$$s_{k} = \underset{s \in \mathcal{V}_{k}}{\operatorname{argmin}} \left\{ \boldsymbol{g}_{k}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \boldsymbol{H}_{k} \boldsymbol{s} + \frac{M}{6} \| \boldsymbol{s} \|^{3} \right\}$$

$$x_{k+1} = \boldsymbol{x}_{k} + \boldsymbol{s}_{k}$$

▶ Equivalently, let $\mathbf{V}_k \in \mathbb{R}^{d imes m}$ whose columns form an orthonormal basis for \mathcal{V}_k

Subspace CRN:
$$z_{k} = \underset{z \in \mathbb{R}^{m}}{\operatorname{argmin}} \left\{ \frac{\tilde{\boldsymbol{g}}_{k}^{\top} \boldsymbol{z} + \frac{1}{2} \boldsymbol{z}^{\top} \tilde{\boldsymbol{H}}_{k} \boldsymbol{z} + \frac{M}{6} \|\boldsymbol{z}\|^{3} \right\}, \quad \boldsymbol{s}_{k} = \boldsymbol{V}_{k} \boldsymbol{z}_{k},$$
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_{k} + \boldsymbol{s}_{k}$$

where $\tilde{\boldsymbol{g}}_k = \boldsymbol{\mathsf{V}}_k^\top \boldsymbol{g}_k \in \mathbb{R}^m$ and $\tilde{\boldsymbol{\mathsf{H}}}_k = \boldsymbol{\mathsf{V}}_k^\top \boldsymbol{\mathsf{H}}_k \boldsymbol{\mathsf{V}}_k \in \mathbb{R}^{m \times m}$

Krylov Cubic Regularized Newton

• How to choose the subspace \mathcal{V}_k ?

- How to choose the subspace \mathcal{V}_k ?
- The existing works choose a random subspace
 - Randomized Block Cubic Newton (RBCN) [Doikov-Richtárik'18]: Sampling a random block of *m* coordinates
 - Stochastic Subspace Cubic Newton (SSCN) [Hanzely et al.'20]: A general random subspace satisfying $\mathbb{E}[\mathbf{V}_k \mathbf{V}_k^{\top}] = \frac{m}{d}\mathbf{I}$

- How to choose the subspace \mathcal{V}_k ?
- The existing works choose a random subspace
 - Randomized Block Cubic Newton (RBCN) [Doikov-Richtárik'18]: Sampling a random block of *m* coordinates
 - Stochastic Subspace Cubic Newton (SSCN) [Hanzely et al.'20]: A general random subspace satisfying $\mathbb{E}[\mathbf{V}_k \mathbf{V}_k^{\top}] = \frac{m}{d}\mathbf{I}$
- Reduced computational costs . . .
 - Computing the subspace Hessian \tilde{H}_k : $\mathcal{O}(m^2)$ in some special cases
 - Solving the cubic subproblem: $\mathcal{O}(m^3)$

- How to choose the subspace \mathcal{V}_k ?
- The existing works choose a random subspace
 - Randomized Block Cubic Newton (RBCN) [Doikov-Richtárik'18]: Sampling a random block of *m* coordinates
 - Stochastic Subspace Cubic Newton (SSCN) [Hanzely et al.'20]: A general random subspace satisfying $\mathbb{E}[\mathbf{V}_k \mathbf{V}_k^{\top}] = \frac{m}{d}\mathbf{I}$
- Reduced computational costs . . .
 - Computing the subspace Hessian \tilde{H}_k : $\mathcal{O}(m^2)$ in some special cases
 - Solving the cubic subproblem: $\mathcal{O}(m^3)$
- ... but a much slower convergence rate

$$\mathcal{O}\left(\frac{d-m}{m}\cdot\frac{1}{k}+\left(\frac{d}{m}\right)^2\cdot\frac{1}{k^2}\right)$$

- How to choose the subspace \mathcal{V}_k ?
- The existing works choose a random subspace
 - Randomized Block Cubic Newton (RBCN) [Doikov-Richtárik'18]: Sampling a random block of *m* coordinates
 - Stochastic Subspace Cubic Newton (SSCN) [Hanzely et al.'20]: A general random subspace satisfying $\mathbb{E}[\mathbf{V}_k \mathbf{V}_k^{\top}] = \frac{m}{d}\mathbf{I}$
- Reduced computational costs . . .
 - Computing the subspace Hessian \tilde{H}_k : $\mathcal{O}(m^2)$ in some special cases
 - Solving the cubic subproblem: $\mathcal{O}(m^3)$
- ... but a much slower convergence rate

$$\mathcal{O}\left(\frac{d-m}{m}\cdot\frac{1}{k}+\left(\frac{d}{m}\right)^2\cdot\frac{1}{k^2}\right)$$

Question: Can we improve the dimensional dependence of subspace second-order methods?

Where the Dimensional Dependence Comes from?

Intuitively, it stems from the fact that the subspace is chosen uniformly random, oblivious to the objective function f

Where the Dimensional Dependence Comes from?

Intuitively, it stems from the fact that the subspace is chosen uniformly random, oblivious to the objective function f

Such a random subspace is unlikely to contain a "good" descent direction

Where the Dimensional Dependence Comes from?

Intuitively, it stems from the fact that the subspace is chosen uniformly random, oblivious to the objective function f

Such a random subspace is unlikely to contain a "good" descent direction

 \blacktriangleright It should be better to employ a subspace customized to the local geometry of f

We propose the Krylov CRN method, where we perform the CRN update over the Krylov subspace associated with H_k and g_k

- ▶ We propose the Krylov CRN method, where we perform the CRN update over the Krylov subspace associated with H_k and g_k
- Can be implemented using the Lanczos method, with 1 gradient evaluation and m Hessian-vector products per iteration

- ▶ We propose the Krylov CRN method, where we perform the CRN update over the Krylov subspace associated with H_k and g_k
- Can be implemented using the Lanczos method, with 1 gradient evaluation and m Hessian-vector products per iteration
- ▶ In the convex case, we prove a dimension-free convergence rate

- ▶ We propose the Krylov CRN method, where we perform the CRN update over the Krylov subspace associated with H_k and g_k
- Can be implemented using the Lanczos method, with 1 gradient evaluation and m Hessian-vector products per iteration
- ▶ In the convex case, we prove a dimension-free convergence rate

Methods	Per-iteration cost	Convergence rate
CRN [Nesterov-Polyak'06]	$\mathcal{O}(d^3)$	$\mathcal{O}(rac{1}{k^2})$
SSCN [Hanzely et al.'20]	$\mathcal{O}(m^3)^*$	$\mathcal{O}(rac{d-m}{m}\cdotrac{1}{k}+rac{d^2}{m^2}\cdotrac{1}{k^2})$
Krylov CRN (ours)	$\mathcal{O}(\mathit{md})^{**}$	$\mathcal{O}(rac{1}{mk}+rac{1}{k^2})$

- ▶ We propose the Krylov CRN method, where we perform the CRN update over the Krylov subspace associated with H_k and g_k
- Can be implemented using the Lanczos method, with 1 gradient evaluation and m Hessian-vector products per iteration
- ▶ In the convex case, we prove a dimension-free convergence rate

Methods	Per-iteration cost	Convergence rate
CRN [Nesterov-Polyak'06]	$\mathcal{O}(d^3)$	$\mathcal{O}(rac{1}{k^2})$
SSCN [Hanzely et al.'20]	$\mathcal{O}(m^3)^*$	$\mathcal{O}(rac{d-m}{m}\cdotrac{1}{k}+rac{d^2}{m^2}\cdotrac{1}{k^2})$
Krylov CRN (ours)	$\mathcal{O}(\mathit{md})^{**}$	$\mathcal{O}(rac{1}{mk}+rac{1}{k^2})$

When the Hessian spectrum possesses certain structure, our method can achieve a faster rate

Ruichen Jiang





Introduction

Proposed Method

③ Convergence Analysis

4 Numerical Results

- We assume that:
 - $f : \mathbb{R}^d \to \mathbb{R}$ is convex
 - f is bounded from below and has bounded level-sets
 - The Hessian of f is Lipchitz, i.e., $\|\nabla^2 f(\mathbf{x}) \nabla^2 f(\mathbf{y})\| \le L_2 \|\mathbf{x} \mathbf{y}\|$, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$

- We assume that:
 - $f : \mathbb{R}^d \to \mathbb{R}$ is convex
 - f is bounded from below and has bounded level-sets
 - The Hessian of f is Lipchitz, i.e., $\|\nabla^2 f(\boldsymbol{x}) \nabla^2 f(\boldsymbol{y})\| \le L_2 \|\boldsymbol{x} \boldsymbol{y}\|$, $\forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d$

An important property: upper bound on quadratic approximation error

$$f(\boldsymbol{x}) - \left(f(\boldsymbol{x}_k) + \boldsymbol{g}_k^\top (\boldsymbol{x} - \boldsymbol{x}_k) + \frac{1}{2} (\boldsymbol{x} - \boldsymbol{x}_k)^\top \boldsymbol{\mathsf{H}}_k (\boldsymbol{x} - \boldsymbol{x}_k)\right) \right| \leq \frac{L_2}{6} \|\boldsymbol{x} - \boldsymbol{x}_k\|^3, \quad \forall \boldsymbol{x} \in \mathbb{R}^d$$

- We assume that:
 - $f : \mathbb{R}^d \to \mathbb{R}$ is convex
 - f is bounded from below and has bounded level-sets
 - The Hessian of f is Lipchitz, i.e., $\|\nabla^2 f(\boldsymbol{x}) \nabla^2 f(\boldsymbol{y})\| \le L_2 \|\boldsymbol{x} \boldsymbol{y}\|$, $\forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d$

An important property: upper bound on quadratic approximation error

$$f(\boldsymbol{x}) - \left(f(\boldsymbol{x}_k) + \boldsymbol{g}_k^\top (\boldsymbol{x} - \boldsymbol{x}_k) + \frac{1}{2} (\boldsymbol{x} - \boldsymbol{x}_k)^\top \boldsymbol{\mathsf{H}}_k (\boldsymbol{x} - \boldsymbol{x}_k)\right) \right| \leq \frac{L_2}{6} \|\boldsymbol{x} - \boldsymbol{x}_k\|^3, \quad \forall \boldsymbol{x} \in \mathbb{R}^d$$

▶ The CRN selects x_{k+1} as the minimizer of the cubic upper approximation

$$\boldsymbol{x}_{k+1} = \operatorname*{argmin}_{\boldsymbol{x} \in \mathbb{R}^d} \left\{ f(\boldsymbol{x}_k) + \boldsymbol{g}_k^\top (\boldsymbol{x} - \boldsymbol{x}_k) + \frac{1}{2} (\boldsymbol{x} - \boldsymbol{x}_k)^\top \boldsymbol{\mathsf{H}}_k (\boldsymbol{x} - \boldsymbol{x}_k) + \frac{L_2}{6} \|\boldsymbol{x} - \boldsymbol{x}_k\|^3 \right\}$$

Theorem

Let $\{\mathbf{x}_k\}_{k\geq 0}$ be generated by CRN and define $D := \sup\{\|\mathbf{x} - \mathbf{x}^*\| : \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$. Then we have

$$f(\mathbf{x}_k) - f^* \le \frac{9L_2D^3}{(k+2)(k+1)}$$

Theorem

Let $\{\mathbf{x}_k\}_{k\geq 0}$ be generated by CRN and define $D := \sup\{\|\mathbf{x} - \mathbf{x}^*\| : \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$. Then we have

$$f(oldsymbol{x}_k) - f^* \leq rac{9L_2D^3}{(k+2)(k+1)}$$

• Using the cubic upper bound $(s_k = x_{k+1} - x_k)$:

$$f(oldsymbol{x}_{k+1}) \leq f(oldsymbol{x}_k) + oldsymbol{g}_k^ op oldsymbol{s}_k + rac{1}{2}oldsymbol{s}_k^ op oldsymbol{\mathsf{H}}_k oldsymbol{s}_k + rac{L_2}{6} \|oldsymbol{s}_k\|^3$$

Theorem

Let $\{\mathbf{x}_k\}_{k\geq 0}$ be generated by CRN and define $D := \sup\{\|\mathbf{x} - \mathbf{x}^*\| : \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$. Then we have

$$f(\mathbf{x}_k) - f^* \leq rac{9L_2D^3}{(k+2)(k+1)}$$

• Using the cubic upper bound
$$(s_k = x_{k+1} - x_k)$$
:

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) + \boldsymbol{g}_k^\top \boldsymbol{s}_k + \frac{1}{2} \boldsymbol{s}_k^\top \boldsymbol{\mathsf{H}}_k \boldsymbol{s}_k + \frac{L_2}{6} \|\boldsymbol{s}_k\|^3$$

Since s_k minimizes the cubic function:

$$f(oldsymbol{x}_{k+1}) \leq f(oldsymbol{x}_k) + oldsymbol{g}_k^{ op} oldsymbol{s} + rac{1}{2} oldsymbol{s}^{ op} oldsymbol{\mathsf{H}}_k oldsymbol{s} + rac{L_2}{6} \|oldsymbol{s}\|^3, \quad orall oldsymbol{s} \in \mathbb{R}^d$$

Theorem

Let $\{\mathbf{x}_k\}_{k\geq 0}$ be generated by CRN and define $D := \sup\{\|\mathbf{x} - \mathbf{x}^*\| : \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$. Then we have

$$f(\mathbf{x}_k) - f^* \leq rac{9L_2D^3}{(k+2)(k+1)}$$

• Using the cubic upper bound
$$(\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k)$$
:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \mathbf{g}_k^{ op} \mathbf{s}_k + rac{1}{2} \mathbf{s}_k^{ op} \mathbf{H}_k \mathbf{s}_k + rac{L_2}{6} \|\mathbf{s}_k\|^3$$

Since s_k minimizes the cubic function:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \mathbf{g}_k^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \mathbf{H}_k \mathbf{s} + \frac{L_2}{6} \|\mathbf{s}\|^3, \quad \forall \mathbf{s} \in \mathbb{R}^d$$

Using the cubic lower bound:

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k + \boldsymbol{s}) + \frac{L_2}{3} \|\boldsymbol{s}\|^3, \quad \forall \boldsymbol{s} \in \mathbb{R}^d$$

Ruichen Jiang

Krylov Cubic Regularized Newton

Using the cubic lower bound:

$$f(oldsymbol{x}_{k+1}) \leq f(oldsymbol{x}_k + oldsymbol{s}) + rac{L_2}{3} \|oldsymbol{s}\|^3, \quad orall oldsymbol{s} \in \mathbb{R}^d$$
Review: Cubic Regularized Newton Method (cont.)

Using the cubic lower bound:

$$f(oldsymbol{x}_{k+1}) \leq f(oldsymbol{x}_k + oldsymbol{s}) + rac{L_2}{3} \|oldsymbol{s}\|^3, \quad orall oldsymbol{s} \in \mathbb{R}^d$$

• Choosing $s = \frac{3}{k+3}(x^* - x_k)$ and using convexity of f:

$$f(\mathbf{x}_{k+1}) \leq f\left(\frac{k}{k+3}\mathbf{x}_{k} + \frac{3}{k+3}\mathbf{x}^{*}\right) + \frac{9L_{2}\|\mathbf{x}_{k} - \mathbf{x}^{*}\|^{3}}{(k+3)^{3}} \leq \frac{k}{k+3}f(\mathbf{x}_{k}) + \frac{3}{k+3}f^{*} + \frac{9L_{2}D^{3}}{(k+3)^{3}}$$

Review: Cubic Regularized Newton Method (cont.)

Using the cubic lower bound:

$$f(oldsymbol{x}_{k+1}) \leq f(oldsymbol{x}_k + oldsymbol{s}) + rac{L_2}{3} \|oldsymbol{s}\|^3, \quad orall oldsymbol{s} \in \mathbb{R}^d$$

• Choosing $s = \frac{3}{k+3}(x^* - x_k)$ and using convexity of f:

$$f(\mathbf{x}_{k+1}) \le f\left(\frac{k}{k+3}\mathbf{x}_k + \frac{3}{k+3}\mathbf{x}^*\right) + \frac{9L_2\|\mathbf{x}_k - \mathbf{x}^*\|^3}{(k+3)^3} \le \frac{k}{k+3}f(\mathbf{x}_k) + \frac{3}{k+3}f^* + \frac{9L_2D^3}{(k+3)^3}$$

The Define $A_k = k(k+1)(k+2)$:

$$\begin{aligned} f(\mathbf{x}_{k+1}) - f^* &\leq \frac{k}{k+3} \left(f(\mathbf{x}_k) - f^* \right) + \frac{9L_2 D^3}{(k+3)^3} \\ \Rightarrow \quad A_{k+1}(f(\mathbf{x}_{k+1}) - f^*) &\leq A_k (f(\mathbf{x}_k) - f^*) + 9L_2 D^3 \\ \Rightarrow \quad A_{k+1}(f(\mathbf{x}_{k+1}) - f^*) &\leq 9L_2 D^3 (k+1) \quad \Rightarrow \quad f(\mathbf{x}_{k+1}) - f^* &\leq \frac{9L_2 D^3}{(k+2)(k+3)} \end{aligned}$$

Subspace CRN:

$$oldsymbol{s}_k = \operatorname*{argmin}_{oldsymbol{s}\in\mathcal{V}_k} \left\{oldsymbol{g}_k^{ op}oldsymbol{s} + rac{1}{2}oldsymbol{s}^{ op}oldsymbol{\mathsf{H}}_koldsymbol{s} + rac{L_2}{6}\|oldsymbol{s}\|^3
ight\}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$$

▶ Where the analysis would break when introducing the subspace V_k ?

► Subspace CRN:

$$s_k = \underset{s \in \mathcal{V}_k}{\operatorname{argmin}} \left\{ g_k^\top s + \frac{1}{2} s^\top H_k s + \frac{L_2}{6} \|s\|^3 \right\}$$
$$x_{k+1} = x_k + s_k$$

$$f(oldsymbol{x}_{k+1}) \leq f(oldsymbol{x}_k) + oldsymbol{g}_k^ op oldsymbol{s} + rac{1}{2}oldsymbol{s}_k^ op oldsymbol{\mathsf{H}}_koldsymbol{s} + rac{L_2}{6}\|oldsymbol{s}\|^3, \hspace{1em} orall oldsymbol{s} \in \mathbb{R}^d$$

Subspace CRN:

$$s_{k} = \underset{s \in \mathcal{V}_{k}}{\operatorname{argmin}} \left\{ \boldsymbol{g}_{k}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \boldsymbol{\mathsf{H}}_{k} \boldsymbol{s} + \frac{L_{2}}{6} \|\boldsymbol{s}\|^{3} \right\}$$
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_{k} + \boldsymbol{s}_{k}$$

• Where the analysis would break when introducing the subspace V_k ?

Recall the crucial inequality

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) + \boldsymbol{g}_k^\top \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}_k^\top \boldsymbol{\mathsf{H}}_k \boldsymbol{s} + \frac{L_2}{6} \|\boldsymbol{s}\|^3, \quad \forall \boldsymbol{s} \in \mathcal{V}_k$$

Subspace CRN:

$$s_k = \underset{s \in \mathcal{V}_k}{\operatorname{argmin}} \left\{ \boldsymbol{g}_k^\top \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^\top \boldsymbol{\mathsf{H}}_k \boldsymbol{s} + \frac{L_2}{6} \|\boldsymbol{s}\|^3 \right\}$$
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{s}_k$$

Where the analysis would break when introducing the subspace V_k?
 Recall the crucial inequality

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) + \boldsymbol{g}_k^\top \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}_k^\top \boldsymbol{\mathsf{H}}_k \boldsymbol{s} + \frac{L_2}{6} \|\boldsymbol{s}\|^3, \quad \forall \boldsymbol{s} \in \mathcal{V}_k$$

▶ Let $\mathbf{P}_k := \mathbf{V}_k \mathbf{V}_k^{\top}$ be the orthogonal projection matrix of \mathcal{V}_k . Since $\mathbf{P}_k \mathbf{s} \in \mathcal{V}_k$, $\forall \mathbf{s} \in \mathbb{R}^d$:

$$f(oldsymbol{x}_{k+1}) \leq f(oldsymbol{x}_k) + oldsymbol{g}_k^ op oldsymbol{P}_k oldsymbol{s} + rac{1}{2} (oldsymbol{P}_k oldsymbol{s})^ op oldsymbol{\mathsf{H}}_k (oldsymbol{P}_k oldsymbol{s}) + rac{L_2}{6} \|oldsymbol{P}_k oldsymbol{s}\|^3, \quad orall oldsymbol{s} \in \mathbb{R}^d$$

- ▶ The exactly same analysis would apply if, for any $s \in \mathbb{R}^d$,
 - (A) $\boldsymbol{g}_k^\top \boldsymbol{\mathsf{P}}_k \boldsymbol{s} \leq \boldsymbol{g}_k^\top \boldsymbol{s}$ (B) $(\boldsymbol{\mathsf{P}}_k \boldsymbol{s})^\top \boldsymbol{\mathsf{H}}_k (\boldsymbol{\mathsf{P}}_k \boldsymbol{s}) \leq \boldsymbol{s}^\top \boldsymbol{\mathsf{H}}_k \boldsymbol{s}$ (C) $\|\boldsymbol{\mathsf{P}}_k \boldsymbol{s}\| \leq \|\boldsymbol{s}\|$

▶ The exactly same analysis would apply if, for any $s \in \mathbb{R}^d$,

(A)
$$\mathbf{g}_{k}^{\top} \mathbf{P}_{k} \mathbf{s} \leq \mathbf{g}_{k}^{\top} \mathbf{s} \iff \mathbf{g}_{k} = \mathbf{P}_{k} \mathbf{g}_{k} \iff \mathbf{g}_{k} \in \mathcal{V}_{k}$$

(B) $(\mathbf{P}_{k} \mathbf{s})^{\top} \mathbf{H}_{k} (\mathbf{P}_{k} \mathbf{s}) \leq \mathbf{s}^{\top} \mathbf{H}_{k} \mathbf{s}$
(C) $\|\mathbf{P}_{k} \mathbf{s}\| \leq \|\mathbf{s}\|$

▶ The exactly same analysis would apply if, for any $s \in \mathbb{R}^d$,

(A)
$$\mathbf{g}_{k}^{\top} \mathbf{P}_{k} \mathbf{s} \leq \mathbf{g}_{k}^{\top} \mathbf{s} \iff \mathbf{g}_{k} = \mathbf{P}_{k} \mathbf{g}_{k} \iff \mathbf{g}_{k} \in \mathcal{V}_{k}$$

(B) $(\mathbf{P}_{k} \mathbf{s})^{\top} \mathbf{H}_{k} (\mathbf{P}_{k} \mathbf{s}) \leq \mathbf{s}^{\top} \mathbf{H}_{k} \mathbf{s} \iff \mathbf{P}_{k} \mathbf{H}_{k} \mathbf{P}_{k} \leq \mathbf{H}_{k} \iff \mathbf{H}_{k} \mathbf{v} \in \mathcal{V}_{k}, \forall \mathbf{v} \in \mathcal{V}_{k}$
(C) $\|\mathbf{P}_{k} \mathbf{s}\| \leq \|\mathbf{s}\|$

▶ The exactly same analysis would apply if, for any $s \in \mathbb{R}^d$,

$$(\mathsf{A}) \ \boldsymbol{g}_k^\top \boldsymbol{\mathsf{P}}_k \boldsymbol{s} \leq \boldsymbol{g}_k^\top \boldsymbol{s} \quad \Leftrightarrow \quad \boldsymbol{g}_k = \boldsymbol{\mathsf{P}}_k \boldsymbol{g}_k \quad \Leftrightarrow \quad \boldsymbol{g}_k \in \mathcal{V}_k$$

(B) $(\mathbf{P}_k s)^\top \mathbf{H}_k (\mathbf{P}_k s) \leq s^\top \mathbf{H}_k s \iff \mathbf{P}_k \mathbf{H}_k \mathbf{P}_k \preceq \mathbf{H}_k \iff \mathbf{H}_k \mathbf{v} \in \mathcal{V}_k, \forall \mathbf{v} \in \mathcal{V}_k$

(C) $\|\mathbf{P}_k \mathbf{s}\| \le \|\mathbf{s}\|$ \checkmark since \mathbf{P}_k is an orthogonal projection matrix

 \blacktriangleright The exactly same analysis would apply if, for any $oldsymbol{s} \in \mathbb{R}^d$,

(A)
$$\boldsymbol{g}_k^\top \boldsymbol{\mathsf{P}}_k \boldsymbol{s} \leq \boldsymbol{g}_k^\top \boldsymbol{s} \quad \Leftrightarrow \quad \boldsymbol{g}_k = \boldsymbol{\mathsf{P}}_k \boldsymbol{g}_k \quad \Leftrightarrow \quad \boldsymbol{g}_k \in \mathcal{V}_k$$

(B) $(\mathbf{P}_k s)^\top \mathbf{H}_k (\mathbf{P}_k s) \leq s^\top \mathbf{H}_k s \iff \mathbf{P}_k \mathbf{H}_k \mathbf{P}_k \preceq \mathbf{H}_k \iff \mathbf{H}_k \mathbf{v} \in \mathcal{V}_k, \forall \mathbf{v} \in \mathcal{V}_k$

(C) $\|\mathbf{P}_k \mathbf{s}\| \le \|\mathbf{s}\|$ \checkmark since \mathbf{P}_k is an orthogonal projection matrix

Combining (A) and (B), we obtain that

 $\boldsymbol{g}_k \in \mathcal{V}_k, \ \boldsymbol{\mathsf{H}}_k \boldsymbol{g}_k \in \mathcal{V}_k, \ \boldsymbol{\mathsf{H}}_k^2 \boldsymbol{g}_k \in \mathcal{V}_k, \ \ldots, \quad \Rightarrow \quad \operatorname{span}\{\boldsymbol{\mathsf{H}}_k^j \boldsymbol{g}_k \mid i = 0, 1, \ldots\} \subset \mathcal{V}_k$

▶ The exactly same analysis would apply if, for any $s \in \mathbb{R}^d$,

$$(\mathsf{A}) \ \mathbf{g}_k^\top \mathbf{P}_k \mathbf{s} \leq \mathbf{g}_k^\top \mathbf{s} \quad \Leftrightarrow \quad \mathbf{g}_k = \mathbf{P}_k \mathbf{g}_k \quad \Leftrightarrow \quad \mathbf{g}_k \in \mathcal{V}_k$$

(B) $(\mathbf{P}_k s)^\top \mathbf{H}_k (\mathbf{P}_k s) \leq s^\top \mathbf{H}_k s \iff \mathbf{P}_k \mathbf{H}_k \mathbf{P}_k \preceq \mathbf{H}_k \iff \mathbf{H}_k \mathbf{v} \in \mathcal{V}_k, \forall \mathbf{v} \in \mathcal{V}_k$

(C) $\|\mathbf{P}_k \mathbf{s}\| \le \|\mathbf{s}\|$ \checkmark since \mathbf{P}_k is an orthogonal projection matrix

Combining (A) and (B), we obtain that

 $\boldsymbol{g}_k \in \mathcal{V}_k, \ \boldsymbol{\mathsf{H}}_k \boldsymbol{g}_k \in \mathcal{V}_k, \ \boldsymbol{\mathsf{H}}_k^2 \boldsymbol{g}_k \in \mathcal{V}_k, \ \dots, \quad \Rightarrow \quad \operatorname{span}\{\boldsymbol{\mathsf{H}}_k^i \boldsymbol{g}_k \mid i = 0, 1, \dots\} \subset \mathcal{V}_k$

▶ This is exactly the maximal Krylov subspace generated by \mathbf{H}_k and \mathbf{g}_k !

▶ Formally, the *j*-th Krylov subspace generated by $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^{d}$ is defined as

$$\mathcal{K}_j(\mathbf{A}, \boldsymbol{b}) = \operatorname{span}\{\boldsymbol{b}, \mathbf{A}\boldsymbol{b}, \dots, \mathbf{A}^{j-1}\boldsymbol{b}\}$$

▶ Formally, the *j*-th Krylov subspace generated by $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^{d}$ is defined as

$$\mathcal{K}_j(\mathbf{A}, \boldsymbol{b}) = \operatorname{span}\{\boldsymbol{b}, \mathbf{A}\boldsymbol{b}, \dots, \mathbf{A}^{j-1}\boldsymbol{b}\}.$$

• Moreover, there exists an integer $r_0 \leq d$ such that:

$$\mathcal{K}_1(\mathbf{A}, \boldsymbol{b}) \subset \mathcal{K}_2(\mathbf{A}, \boldsymbol{b}) \subset \cdots \subset \underbrace{\mathcal{K}_{r_0}(\mathbf{A}, \boldsymbol{b})}_{\dim = r_0} = \mathcal{K}_{r_0+1}(\mathbf{A}, \boldsymbol{b}) = \mathcal{K}_{r_0+2}(\mathbf{A}, \boldsymbol{b}) = \cdots$$

We call $\mathcal{K}_{r_0}(\mathbf{A}, \mathbf{b})$ the maximal Krylov subspace

Formally, the *j*-th Krylov subspace generated by $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^{d}$ is defined as

$$\mathcal{K}_j(\mathbf{A}, \boldsymbol{b}) = \operatorname{span}\{\boldsymbol{b}, \mathbf{A}\boldsymbol{b}, \dots, \mathbf{A}^{j-1}\boldsymbol{b}\}.$$

• Moreover, there exists an integer $r_0 \leq d$ such that:

$$\mathcal{K}_1(\mathbf{A}, \boldsymbol{b}) \subset \mathcal{K}_2(\mathbf{A}, \boldsymbol{b}) \subset \cdots \subset \underbrace{\mathcal{K}_{r_0}(\mathbf{A}, \boldsymbol{b})}_{\dim = r_0} = \mathcal{K}_{r_0+1}(\mathbf{A}, \boldsymbol{b}) = \mathcal{K}_{r_0+2}(\mathbf{A}, \boldsymbol{b}) = \cdots$$

We call $\mathcal{K}_{r_0}(\mathbf{A}, \mathbf{b})$ the maximal Krylov subspace

▶ To sum up: if we let $V_k = K_{r_0}(\mathbf{H}_k, \mathbf{g}_k)$, the subspace CRN retains the same convergence rate of CRN

Formally, the *j*-th Krylov subspace generated by $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^{d}$ is defined as

$$\mathcal{K}_j(\mathbf{A}, \mathbf{b}) = \operatorname{span}\{\mathbf{b}, \mathbf{A}\mathbf{b}, \dots, \mathbf{A}^{j-1}\mathbf{b}\}.$$

• Moreover, there exists an integer $r_0 \leq d$ such that:

$$\mathcal{K}_1(\mathbf{A}, \boldsymbol{b}) \subset \mathcal{K}_2(\mathbf{A}, \boldsymbol{b}) \subset \cdots \subset \underbrace{\mathcal{K}_{r_0}(\mathbf{A}, \boldsymbol{b})}_{\dim = r_0} = \mathcal{K}_{r_0+1}(\mathbf{A}, \boldsymbol{b}) = \mathcal{K}_{r_0+2}(\mathbf{A}, \boldsymbol{b}) = \cdots$$

We call $\mathcal{K}_{r_0}(\mathbf{A}, \mathbf{b})$ the maximal Krylov subspace

- ▶ To sum up: if we let $V_k = K_{r_0}(\mathbf{H}_k, \mathbf{g}_k)$, the subspace CRN retains the same convergence rate of CRN
- However, r_0 can be as large as $d \Rightarrow$ we use the Krylov subspace up to dim m

Let us check if V_k = K_m(H_k, g_k) satisfies the three conditions:
(A) g_k^TP_ks ≤ g_k^Ts ⇔ g_k ∈ V_k
(B) (P_ks)^TH_k(P_ks) ≤ s^TH_ks ⇔ H_kv ∈ V_k, ∀v ∈ V_k
(C) ||P_ks|| ≤ ||s||

Let us check if V_k = K_m(H_k, g_k) satisfies the three conditions:
(A) g_k^TP_ks ≤ g_k^Ts ⇔ g_k ∈ V_k ✓
(B) (P_ks)^TH_k(P_ks) ≤ s^TH_ks ⇔ H_kv ∈ V_k, ∀v ∈ V_k
(C) ||P_ks|| ≤ ||s||

Let us check if V_k = K_m(H_k, g_k) satisfies the three conditions:
(A) g_k^TP_ks ≤ g_k^Ts ⇔ g_k ∈ V_k ✓
(B) (P_ks)^TH_k(P_ks) ≤ s^TH_ks ⇔ H_kv ∈ V_k, ∀v ∈ V_k
(C) ||P_ks|| ≤ ||s|| ✓ since P_k is an orthogonal projection matrix

Let us check if V_k = K_m(H_k, g_k) satisfies the three conditions:
(A) g_k^TP_ks ≤ g_k^Ts ⇔ g_k ∈ V_k ✓
(B) (P_ks)^TH_k(P_ks) ≤ s^TH_ks ⇔ H_kv ∈ V_k, ∀v ∈ V_k × H_k^{m-1}g_k ∈ V_k but H_k^mg_k ∉ V_k
(C) ||P_ks|| ≤ ||s|| ✓ since P_k is an orthogonal projection matrix

• Let us check if $\mathcal{V}_k = \mathcal{K}_m(\mathbf{H}_k, \mathbf{g}_k)$ satisfies the three conditions:

 $(\mathsf{A}) \ \boldsymbol{g}_k^\top \mathsf{P}_k \boldsymbol{s} \leq \boldsymbol{g}_k^\top \boldsymbol{s} \quad \Leftrightarrow \quad \boldsymbol{g}_k \in \mathcal{V}_k \quad \checkmark$

(B) $(\mathbf{P}_k s)^\top \mathbf{H}_k (\mathbf{P}_k s) \leq s^\top \mathbf{H}_k s \iff \mathbf{H}_k v \in \mathcal{V}_k, \forall v \in \mathcal{V}_k \quad \mathbf{X} \ \mathbf{H}_k^{m-1} \mathbf{g}_k \in \mathcal{V}_k \text{ but } \mathbf{H}_k^m \mathbf{g}_k \notin \mathcal{V}_k$

(C) $\|\mathbf{P}_k \mathbf{s}\| \le \|\mathbf{s}\|$ \checkmark since \mathbf{P}_k is an orthogonal projection matrix

Only need to control the error from Condition (B)

• Let us check if $\mathcal{V}_k = \mathcal{K}_m(\mathbf{H}_k, \mathbf{g}_k)$ satisfies the three conditions:

$$\textbf{(A)} \ \boldsymbol{g}_k^\top \boldsymbol{\mathsf{P}}_k \boldsymbol{s} \leq \boldsymbol{g}_k^\top \boldsymbol{s} \quad \Leftrightarrow \quad \boldsymbol{g}_k \in \mathcal{V}_k \quad \checkmark$$

(B) $(\mathbf{P}_k s)^\top \mathbf{H}_k (\mathbf{P}_k s) \leq s^\top \mathbf{H}_k s \iff \mathbf{H}_k v \in \mathcal{V}_k, \forall v \in \mathcal{V}_k \quad \mathbf{X} \ \mathbf{H}_k^{m-1} \mathbf{g}_k \in \mathcal{V}_k \text{ but } \mathbf{H}_k^m \mathbf{g}_k \notin \mathcal{V}_k$

(C) $\|\mathbf{P}_k \mathbf{s}\| \le \|\mathbf{s}\|$ \checkmark since \mathbf{P}_k is an orthogonal projection matrix

- Only need to control the error from Condition (B)
- ► In comparison, when V_k is chosen randomly, both (A) and (B) only hold approximately and the induced errors depend on m/d

- ▶ Input: Initial point $\mathbf{x}_0 \in \mathbb{R}^d$, subspace dimension m, regularization parameter M > 0
- for k = 0, 1, ..., do
 - $\mathbf{V}_k \leftarrow$ the orthnormal basis of $\mathcal{K}_m(\mathbf{H}_k, \mathbf{g}_k)$, $\tilde{\mathbf{g}}_k \leftarrow \mathbf{V}_k^\top \mathbf{g}_k$, $\tilde{\mathbf{H}}_k \leftarrow \mathbf{V}_k^\top \mathbf{H}_k \mathbf{V}_k$
 - Solve the cubic subproblem

$$oldsymbol{z}_k = \operatorname*{argmin}_{oldsymbol{z} \in \mathbb{R}^m} \Big\{ \widetilde{oldsymbol{g}}_k^{ op} oldsymbol{z} + rac{1}{2} oldsymbol{z}^{ op} \widetilde{oldsymbol{H}}_k oldsymbol{z} + rac{M}{6} \|oldsymbol{z}\|^3 \Big\},$$

• Update
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{V}_k \boldsymbol{z}_k$$

end for

- ▶ Input: Initial point $x_0 \in \mathbb{R}^d$, subspace dimension m, regularization parameter M > 0
- for k = 0, 1, ..., do
 - $\mathbf{V}_k \leftarrow$ the orthnormal basis of $\mathcal{K}_m(\mathbf{H}_k, \mathbf{g}_k)$, $\tilde{\mathbf{g}}_k \leftarrow \mathbf{V}_k^\top \mathbf{g}_k$, $\tilde{\mathbf{H}}_k \leftarrow \mathbf{V}_k^\top \mathbf{H}_k \mathbf{V}_k$
 - Solve the cubic subproblem

$$oldsymbol{z}_k = \operatorname*{argmin}_{oldsymbol{z} \in \mathbb{R}^m} \Big\{ \widetilde{oldsymbol{g}}_k^{ op} oldsymbol{z} + rac{1}{2} oldsymbol{z}^{ op} \widetilde{oldsymbol{H}}_k oldsymbol{z} + rac{M}{6} \|oldsymbol{z}\|^3 \Big\},$$

• Update
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{\mathsf{V}}_k \boldsymbol{z}_k$$

end for

▶ The remaining question: How to compute \mathbf{V}_k , $\tilde{\mathbf{g}}_k$, and $\tilde{\mathbf{H}}_k$?

- ▶ Input: Initial point $\mathbf{x}_0 \in \mathbb{R}^d$, subspace dimension m, regularization parameter M > 0
- for k = 0, 1, ..., do
 - $\mathbf{V}_k \leftarrow$ the orthnormal basis of $\mathcal{K}_m(\mathbf{H}_k, \mathbf{g}_k)$, $\tilde{\mathbf{g}}_k \leftarrow \mathbf{V}_k^\top \mathbf{g}_k$, $\tilde{\mathbf{H}}_k \leftarrow \mathbf{V}_k^\top \mathbf{H}_k \mathbf{V}_k$
 - Solve the cubic subproblem

$$oldsymbol{z}_k = \operatorname*{argmin}_{oldsymbol{z} \in \mathbb{R}^m} \Big\{ \widetilde{oldsymbol{g}}_k^{ op} oldsymbol{z} + rac{1}{2} oldsymbol{z}^{ op} \widetilde{oldsymbol{H}}_k oldsymbol{z} + rac{M}{6} \|oldsymbol{z}\|^3 \Big\},$$

• Update
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{V}_k \boldsymbol{z}_k$$

end for

▶ The remaining question: How to compute V_k , \tilde{g}_k , and \tilde{H}_k ? \leftarrow Lanczos method!

Lanczos Method

Lanczos Method

Input: H∈ ℝ^{d×d}, g∈ ℝ^d, and the dimension m
Initialize: v₁ = g/||g||, β₁ = 0, v₀ = 0
for j = 1, 2, ..., m do
w_j ← Hv_j - β_jv_{j-1} // one Hessian-vector product (HVP)
α_j ← w_j^Tv_j
w_j ← w_j - α_jv_j
β_{j+1} ← ||w_j||₂
v_{j+1} ← w_j/β_{j+1}
end for
Output: V = [v₁, v₂, ..., v_m], ğ = ||g||e₁, and H =

$$\begin{bmatrix} α1 β2 β3 ↔ ↔ β3 ↔ ↔ βm-1 βm αm βm αm \end{bmatrix}$$

► No need to compute & store H explicitly; only requires m HVPs ⇒ Can be done efficiently via back-propagation [Pearlmutter'94]

Lanczos Method

Input:
$$\mathbf{H} \in \mathbb{R}^{d \times d}$$
, $\mathbf{g} \in \mathbb{R}^{d}$, and the dimension m
Initialize: $\mathbf{v}_{1} = \mathbf{g}/||\mathbf{g}||$, $\beta_{1} = 0$, $\mathbf{v}_{0} = 0$
for $j = 1, 2, \ldots, m$ do
 $\mathbf{w}_{j} \leftarrow \mathbf{H}\mathbf{v}_{j} - \beta_{j}\mathbf{v}_{j-1}$ // one Hessian-vector product (HVP)
 $\alpha_{j} \leftarrow \mathbf{w}_{j}^{\top}\mathbf{v}_{j}$
 $\mathbf{w}_{j} \leftarrow \mathbf{w}_{j} - \alpha_{j}\mathbf{v}_{j}$
 $\beta_{j+1} \leftarrow ||\mathbf{w}_{j}||_{2}$
 $\mathbf{v}_{j+1} \leftarrow \mathbf{w}_{j}/\beta_{j+1}$
end for
Output: $\mathbf{V} = [\mathbf{v}_{1}, \mathbf{v}_{2}, \ldots, \mathbf{v}_{m}]$, $\tilde{\mathbf{g}} = ||\mathbf{g}||\mathbf{e}_{1}$, and $\tilde{\mathbf{H}} = \begin{bmatrix} \alpha_{1} & \beta_{2} & \beta_{3} & \beta_{3} & \ddots & \ddots & \beta_{m-1} & \beta_{m-1} & \alpha_{m-1} & \beta_{m} & \beta_{m-1} & \alpha_{m-1} & \beta_{m} & \beta_{m} & \alpha_{m} \end{bmatrix}$

▶ No need to compute & store **H** explicitly; only requires *m* HVPs

 \Rightarrow Can be done efficiently via back-propagation [Pearlmutter'94]

Bonus: H
 is a sparse tridiagonal matrix

Ruichen Jiang

Krylov Cubic Regularized Newton

- **Input:** Initial point $\mathbf{x}_0 \in \mathbb{R}^d$, subspace dimension m, regularization parameter M > 0
- for k = 0, 1, ..., do
 - $(\mathbf{V}_k, \tilde{\mathbf{g}}_k, \tilde{\mathbf{H}}_k) \leftarrow \text{Lanczos}(\mathbf{H}_k, \mathbf{g}_k; m)$
 - Solve the cubic subproblem

$$oldsymbol{z}_k = \operatorname*{argmin}_{oldsymbol{z} \in \mathbb{R}^m} \Big\{ oldsymbol{ ilde{g}}_k^{ op} oldsymbol{z} + rac{1}{2} oldsymbol{z}^{ op} oldsymbol{ ilde{H}}_k oldsymbol{z} + rac{M}{6} \|oldsymbol{z}\|^3 \Big\},$$

• Update
$$oldsymbol{x}_{k+1} = oldsymbol{x}_k + oldsymbol{\mathsf{V}}_k oldsymbol{z}_k$$

end for

- Per-iteration computational cost
 - Performing Lanczos iteration: $m \text{ HVPs} \Rightarrow \mathcal{O}(md)$
 - Solving the cubic subproblem: $\mathcal{O}(m)$



Introduction

② The Proposed Method

③ Convergence Analysis

4 Numerical Results

▶ For simplicity, we further assume that $\nabla^2 f(\mathbf{x}) \preceq L_1 \mathbf{I}$ as in [Hanzely et al.'20]

▶ For simplicity, we further assume that $\nabla^2 f(\mathbf{x}) \preceq L_1 \mathbf{I}$ as in [Hanzely et al.'20]

Lemma

Let $\{\mathbf{x}_k\}$ be generated by Krylov CRN with subspace dim m. We have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \mathbf{g}_k^{\top} \mathbf{s} + \frac{1}{2} \mathbf{s}^{\top} \mathbf{H}_k \mathbf{s} + \frac{L_2}{6} \|\mathbf{s}\|^3 + \frac{L_1}{2m} \|\mathbf{s}\|^2, \quad \forall \mathbf{s} \in \mathbb{R}^d.$$

▶ For simplicity, we further assume that $\nabla^2 f(\mathbf{x}) \leq L_1 \mathbf{I}$ as in [Hanzely et al.'20]

Lemma

Let $\{\mathbf{x}_k\}$ be generated by Krylov CRN with subspace dim m. We have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \mathbf{g}_k^{\top} \mathbf{s} + \frac{1}{2} \mathbf{s}^{\top} \mathbf{H}_k \mathbf{s} + \frac{L_2}{6} \|\mathbf{s}\|^3 + \frac{L_1}{2m} \|\mathbf{s}\|^2, \quad \forall \mathbf{s} \in \mathbb{R}^d.$$

The additional error term is independent of d

▶ For simplicity, we further assume that $\nabla^2 f(\mathbf{x}) \leq L_1 \mathbf{I}$ as in [Hanzely et al.'20]

Lemma

Let $\{\mathbf{x}_k\}$ be generated by Krylov CRN with subspace dim m. We have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \mathbf{g}_k^{\top} \mathbf{s} + \frac{1}{2} \mathbf{s}^{\top} \mathbf{H}_k \mathbf{s} + \frac{L_2}{6} \|\mathbf{s}\|^3 + \frac{L_1}{2m} \|\mathbf{s}\|^2, \quad \forall \mathbf{s} \in \mathbb{R}^d.$$

- The additional error term is independent of d
- It diminishes as m increases

Proof Sketch

- For simplicity, ignore the subscript k
- ▶ Let $\mathbf{P}^{(m)}$ be the orthogonal projection matrix associated with $\mathcal{K}_m(\mathbf{H}, \boldsymbol{g})$

Proof Sketch

- For simplicity, ignore the subscript k
- ▶ Let $\mathbf{P}^{(m)}$ be the orthogonal projection matrix associated with $\mathcal{K}_m(\mathbf{H}, \boldsymbol{g})$
- Recall the key inequality: for any $\boldsymbol{s} \in \mathbb{R}^d$,

$$f(\boldsymbol{x}_{+}) \leq f(\boldsymbol{x}) + \boldsymbol{g}^{\top} \boldsymbol{\mathsf{P}}^{(m)} \boldsymbol{s} + \frac{1}{2} (\boldsymbol{\mathsf{P}}^{(m)} \boldsymbol{s})^{\top} \boldsymbol{\mathsf{H}} (\boldsymbol{\mathsf{P}}^{(m)} \boldsymbol{s}) + \frac{L_2}{6} \| \boldsymbol{\mathsf{P}}^{(m)} \boldsymbol{s} \|^3$$
Proof Sketch

- For simplicity, ignore the subscript k
- ▶ Let $\mathbf{P}^{(m)}$ be the orthogonal projection matrix associated with $\mathcal{K}_m(\mathbf{H}, \boldsymbol{g})$
- Recall the key inequality: for any $\boldsymbol{s} \in \mathbb{R}^d$,

$$f(\mathbf{x}_+) \leq f(\mathbf{x}) + \mathbf{g}^{ op} \mathbf{P}^{(m)} \mathbf{s} + rac{1}{2} (\mathbf{P}^{(m)} \mathbf{s})^{ op} \mathbf{H} (\mathbf{P}^{(m)} \mathbf{s}) + rac{L_2}{6} \|\mathbf{P}^{(m)} \mathbf{s}\|^3$$

▶ Since $g \in \mathcal{K}_m(\mathbf{H}, g) \iff \mathbf{P}^{(m)}g = g$ and $\|\mathbf{P}^{(m)}s\| \le \|s\|$:

$$f(\boldsymbol{x}_{+}) \leq f(\boldsymbol{x}) + \boldsymbol{g}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \boldsymbol{\mathsf{H}} \boldsymbol{s} + \frac{L_2}{6} \|\boldsymbol{s}\|^3 + \frac{1}{2} \boldsymbol{s}^{\top} \left(\boldsymbol{\mathsf{P}}^{(m)} \boldsymbol{\mathsf{H}} \boldsymbol{\mathsf{P}}^{(m)} - \boldsymbol{\mathsf{H}} \right) \boldsymbol{s}$$

Proof Sketch

- For simplicity, ignore the subscript k
- ▶ Let $\mathbf{P}^{(m)}$ be the orthogonal projection matrix associated with $\mathcal{K}_m(\mathbf{H}, \boldsymbol{g})$
- Recall the key inequality: for any $\boldsymbol{s} \in \mathbb{R}^d$,

$$f(\mathbf{x}_+) \leq f(\mathbf{x}) + \mathbf{g}^{ op} \mathbf{P}^{(m)} \mathbf{s} + rac{1}{2} (\mathbf{P}^{(m)} \mathbf{s})^{ op} \mathbf{H} (\mathbf{P}^{(m)} \mathbf{s}) + rac{L_2}{6} \|\mathbf{P}^{(m)} \mathbf{s}\|^3$$

▶ Since $m{g} \in \mathcal{K}_m(\mathbf{H}, m{g}) \iff \mathbf{P}^{(m)} m{g} = m{g}$ and $\|\mathbf{P}^{(m)} m{s}\| \le \|m{s}\|$:

$$f(\boldsymbol{x}_{+}) \leq f(\boldsymbol{x}) + \boldsymbol{g}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \boldsymbol{\mathsf{H}} \boldsymbol{s} + \frac{L_2}{6} \|\boldsymbol{s}\|^3 + \frac{1}{2} \boldsymbol{s}^{\top} \left(\boldsymbol{\mathsf{P}}^{(m)} \boldsymbol{\mathsf{H}} \boldsymbol{\mathsf{P}}^{(m)} - \boldsymbol{\mathsf{H}} \right) \boldsymbol{s}$$

▶ Recall the Lanczos algorithm generates $\{\mathbf{v}_j\}_{j=1}^m$ and $\{\beta_{j+1}\}_{j=1}^m$:

$$\frac{1}{2}\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(m)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(m)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s}\leq\beta_{m+1}|\boldsymbol{\mathsf{v}}_{m}^{\top}\boldsymbol{s}||\boldsymbol{\mathsf{v}}_{m+1}^{\top}\boldsymbol{s}|$$

▶ It can be shown that $\beta_{j+1} \leq L_1/2$ for all $j \geq 1$:

$$\frac{1}{2}\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(m)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(m)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s} \leq \beta_{m+1}|\boldsymbol{v}_{m}^{\top}\boldsymbol{s}||\boldsymbol{v}_{m+1}^{\top}\boldsymbol{s}| \leq \frac{L_{1}}{2}\|\boldsymbol{s}\|^{2}$$

lt can be shown that
$$\beta_{j+1} \leq L_1/2$$
 for all $j \geq 1$:

$$\frac{1}{2}\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(m)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(m)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s} \leq \beta_{m+1}|\boldsymbol{v}_{m}^{\top}\boldsymbol{s}||\boldsymbol{v}_{m+1}^{\top}\boldsymbol{s}| \leq \frac{L_{1}}{2}\|\boldsymbol{s}\|^{2}$$

 \blacktriangleright However, the error term does not diminish when *m* increases

• It can be shown that $\beta_{j+1} \leq L_1/2$ for all $j \geq 1$:

$$\frac{1}{2}\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(m)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(m)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s} \leq \beta_{m+1}|\boldsymbol{v}_{m}^{\top}\boldsymbol{s}||\boldsymbol{v}_{m+1}^{\top}\boldsymbol{s}| \leq \frac{L_{1}}{2}\|\boldsymbol{s}\|^{2}$$

However, the error term does not diminish when m increases

Turns out we can strengthen the bound to

$$f(\boldsymbol{x}_{+}) \leq f(\boldsymbol{x}) + \boldsymbol{g}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \boldsymbol{\mathsf{H}} \boldsymbol{s} + \frac{L_{2}}{6} \|\boldsymbol{s}\|^{3} + \frac{1}{2} \min_{j \in \{1,...,m\}} \left\{ \boldsymbol{s}^{\top} \left(\boldsymbol{\mathsf{P}}^{(j)} \boldsymbol{\mathsf{H}} \boldsymbol{\mathsf{P}}^{(j)} - \boldsymbol{\mathsf{H}} \right) \boldsymbol{s} \right\},$$

where $\mathbf{P}^{(j)}$ be the orthogonal projection matrix associated with $\mathcal{K}_{i}(\mathbf{H}, \boldsymbol{g})$

• It can be shown that $\beta_{j+1} \leq L_1/2$ for all $j \geq 1$:

$$\frac{1}{2}\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(m)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(m)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s} \leq \beta_{m+1}|\boldsymbol{v}_{m}^{\top}\boldsymbol{s}||\boldsymbol{v}_{m+1}^{\top}\boldsymbol{s}| \leq \frac{L_{1}}{2}\|\boldsymbol{s}\|^{2}$$

However, the error term does not diminish when m increases

Turns out we can strengthen the bound to

$$f(\boldsymbol{x}_{+}) \leq f(\boldsymbol{x}) + \boldsymbol{g}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \boldsymbol{\mathsf{H}} \boldsymbol{s} + \frac{L_{2}}{6} \|\boldsymbol{s}\|^{3} + \frac{1}{2} \min_{j \in \{1,...,m\}} \left\{ \boldsymbol{s}^{\top} \left(\boldsymbol{\mathsf{P}}^{(j)} \boldsymbol{\mathsf{H}} \boldsymbol{\mathsf{P}}^{(j)} - \boldsymbol{\mathsf{H}} \right) \boldsymbol{s} \right\},$$

where $\mathbf{P}^{(j)}$ be the orthogonal projection matrix associated with $\mathcal{K}_j(\mathbf{H}, \mathbf{g})$ \blacktriangleright Hence,

$$\frac{1}{2}\min_{j\in\{1,...,m\}}\left\{\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(j)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(j)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s}\right\}\leq\frac{L_{1}}{2}\min_{j\in\{1,...,m\}}|\boldsymbol{v}_{j}^{\top}\boldsymbol{s}||\boldsymbol{v}_{j+1}^{\top}\boldsymbol{s}|$$

• It can be shown that $\beta_{j+1} \leq L_1/2$ for all $j \geq 1$:

$$\frac{1}{2}\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(m)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(m)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s} \leq \beta_{m+1}|\boldsymbol{v}_{m}^{\top}\boldsymbol{s}||\boldsymbol{v}_{m+1}^{\top}\boldsymbol{s}| \leq \frac{L_{1}}{2}\|\boldsymbol{s}\|^{2}$$

However, the error term does not diminish when m increases

Turns out we can strengthen the bound to

$$f(\mathbf{x}_+) \leq f(\mathbf{x}) + \mathbf{g}^{\top}\mathbf{s} + rac{1}{2}\mathbf{s}^{\top}\mathbf{H}\mathbf{s} + rac{L_2}{6}\|\mathbf{s}\|^3 + rac{1}{2}\min_{j\in\{1,...,m\}}\left\{\mathbf{s}^{\top}\left(\mathbf{P}^{(j)}\mathbf{H}\mathbf{P}^{(j)}-\mathbf{H}
ight)\mathbf{s}
ight\},$$

where $\mathbf{P}^{(j)}$ be the orthogonal projection matrix associated with $\mathcal{K}_j(\mathbf{H}, \boldsymbol{g})$ \blacktriangleright Hence,

$$\frac{1}{2}\min_{j\in\{1,...,m\}}\left\{\boldsymbol{s}^{\top}\left(\boldsymbol{\mathsf{P}}^{(j)}\boldsymbol{\mathsf{H}}\boldsymbol{\mathsf{P}}^{(j)}-\boldsymbol{\mathsf{H}}\right)\boldsymbol{s}\right\}\leq\frac{L_{1}}{2}\min_{j\in\{1,...,m\}}|\boldsymbol{v}_{j}^{\top}\boldsymbol{s}||\boldsymbol{v}_{j+1}^{\top}\boldsymbol{s}|$$

► Since $\{\mathbf{v}_j\}$ are orthonormal, we further have $\min_{j \in \{1,...,m\}} |\mathbf{v}_j^\top \mathbf{s}| |\mathbf{v}_{j+1}^\top \mathbf{s}| \le \frac{1}{m} \|\mathbf{s}\|^2$

Let $\{\mathbf{x}_k\}$ be generated by the Krylov CRN method. Then we have

$$f(\mathbf{x}_k) - f^* \le \frac{9L_1D^2}{2mk} + \frac{9L_2D^3}{k^2}$$

► The convergence rate is independent of *d*

Let $\{\mathbf{x}_k\}$ be generated by the Krylov CRN method. Then we have

$$f(\mathbf{x}_k) - f^* \le \frac{9L_1D^2}{2mk} + \frac{9L_2D^3}{k^2}$$

- ► The convergence rate is independent of *d*
- ► To achieve accuracy ϵ , the number of iterations required is $\mathcal{O}\left(\frac{1}{m\epsilon} + \frac{1}{\sqrt{\epsilon}}\right)$

Let $\{\mathbf{x}_k\}$ be generated by the Krylov CRN method. Then we have $9I_1D^2 = 9I_2D^3$

$$f(\mathbf{x}_k) - f^* \leq \frac{9L_1D^2}{2mk} + \frac{9L_2D^2}{k^2}$$

- The convergence rate is independent of d
- ► To achieve accuracy ϵ , the number of iterations required is $\mathcal{O}\left(\frac{1}{m\epsilon} + \frac{1}{\sqrt{\epsilon}}\right)$
- ▶ In comparison, SSCN [Hanzely et al.'20] requires $\mathcal{O}(\frac{d-m}{m} \cdot \frac{1}{\epsilon} + \frac{d}{m} \cdot \frac{1}{\sqrt{\epsilon}})$

Let $\{\mathbf{x}_k\}$ be generated by the Krylov CRN method. Then we have $f(\mathbf{x}_k) - f^* \leq \frac{9L_1D^2}{2mk} + \frac{9L_2D^3}{k^2}$

- The convergence rate is independent of d
- ► To achieve accuracy ϵ , the number of iterations required is $\mathcal{O}\left(\frac{1}{m\epsilon} + \frac{1}{\sqrt{\epsilon}}\right)$
- ▶ In comparison, SSCN [Hanzely et al.'20] requires $\mathcal{O}(\frac{d-m}{m} \cdot \frac{1}{\epsilon} + \frac{d}{m} \cdot \frac{1}{\sqrt{\epsilon}})$
- When $m \ll d$, our complexity is lower by a factor of d

Main Result: Strongly Convex Setting

▶ We also consider the strongly convex setting, i.e., $\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}$, $\forall \mathbf{x} \in \mathbb{R}^d$

Theorem

Let $\{\mathbf{x}_k\}$ be generated by the Krylov CRN method. Then, the number of iterations required to reach $\delta_k := f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$ can be upper bounded by

$$k = \mathcal{O}igg(igg(rac{L_1}{m\mu}+1igg)\lograc{\delta_0}{\epsilon}+rac{\sqrt{L_2}\delta_0^{0.25}}{\mu^{0.75}}igg)$$

► In comparison, SSCN requires $\mathcal{O}\left(\left(\frac{d-m}{m}\frac{L_1}{\mu} + \frac{d}{m}\right)\log\frac{1}{\epsilon}\right)$

Again, when m le d, we shave a factor of d

 \blacktriangleright The convergence bound depends on L_1 , i.e., the largest eigenvalue of the Hessian

- \blacktriangleright The convergence bound depends on L_1 , i.e., the largest eigenvalue of the Hessian
- With a more refined analysis, we can replace L_1 by

$$\rho_{\max}^{(m)} := \max_{i \in \{0,1,\dots,k-1\}} \{ \rho^{(m)}(\mathbf{H}_i, \mathbf{g}_i) \},$$

- \blacktriangleright The convergence bound depends on L_1 , i.e., the largest eigenvalue of the Hessian
- With a more refined analysis, we can replace L_1 by

$$\rho_{\max}^{(m)} := \max_{i \in \{0,1,\dots,k-1\}} \{ \rho^{(m)}(\mathbf{H}_i, \mathbf{g}_i) \},$$

where $\rho^{(m)}(\mathbf{H}, \mathbf{g})$ is defined by

$$\rho^{(m)}(\mathbf{H}, \mathbf{g}) = \min_{c_0, \dots, c_{m-1} \in \mathbb{R}} \left\| \mathbf{H}^m \frac{\mathbf{g}}{\|\mathbf{g}\|} - \sum_{i=0}^{m-1} c_i \mathbf{H}^i \frac{\mathbf{g}}{\|\mathbf{g}\|} \right\|^{\frac{1}{m}}$$

- The convergence bound depends on L_1 , i.e., the largest eigenvalue of the Hessian
- With a more refined analysis, we can replace L_1 by

$$\rho_{\max}^{(m)} := \max_{i \in \{0,1,\dots,k-1\}} \{ \rho^{(m)}(\mathbf{H}_i, \mathbf{g}_i) \},$$

where $\rho^{(m)}(\mathbf{H}, \mathbf{g})$ is defined by

$$\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) = \min_{c_0, \dots, c_{m-1} \in \mathbb{R}} \left\| \mathbf{H}^m \frac{\boldsymbol{g}}{\|\boldsymbol{g}\|} - \sum_{i=0}^{m-1} c_i \mathbf{H}^i \frac{\boldsymbol{g}}{\|\boldsymbol{g}\|} \right\|^{\frac{1}{m}}$$

For a polynomial $p(x) = x^m - \sum_{i=0}^{m-1} c_i x^i$, define $p(\mathbf{H}) = \mathbf{H}^m - \sum_{i=0}^{m-1} c_i \mathbf{H}^i$

- \blacktriangleright The convergence bound depends on L_1 , i.e., the largest eigenvalue of the Hessian
- With a more refined analysis, we can replace L_1 by

$$\rho_{\max}^{(m)} := \max_{i \in \{0,1,\dots,k-1\}} \{ \rho^{(m)}(\mathbf{H}_i, \mathbf{g}_i) \},\$$

where $\rho^{(m)}(\mathbf{H}, \mathbf{g})$ is defined by

$$\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) = \min_{c_0, \dots, c_{m-1} \in \mathbb{R}} \left\| \mathbf{H}^m \frac{\boldsymbol{g}}{\|\boldsymbol{g}\|} - \sum_{i=0}^{m-1} c_i \mathbf{H}^i \frac{\boldsymbol{g}}{\|\boldsymbol{g}\|} \right\|^{\frac{1}{m}}$$

For a polynomial p(x) = x^m − ∑_{i=0}^{m-1} c_ixⁱ, define p(H) = H^m − ∑_{i=0}^{m-1} c_iHⁱ
 Let M_m be the set of monic polynomials of degree m. Then we have

$$\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) = \min_{\boldsymbol{p} \in \mathcal{M}_m} \left\| \boldsymbol{p}(\mathbf{H}) \frac{\boldsymbol{g}}{\|\boldsymbol{g}\|} \right\|^{\frac{1}{m}} \leq \min_{\boldsymbol{p} \in \mathcal{M}_m} \|\boldsymbol{p}(\mathbf{H})\|^{\frac{1}{m}}$$

1

Assume that the Hessian **H** has *r* distinct eigenvalues in decreasing order: $\lambda_1 > \lambda_2 > \cdots > \lambda_r$

Assume that the Hessian **H** has *r* distinct eigenvalues in decreasing order: $\lambda_1 > \lambda_2 > \cdots > \lambda_r$

Then

$$\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) \leq \min_{\boldsymbol{p} \in \mathcal{M}_m} \|\boldsymbol{p}(\mathbf{H})\|^{\frac{1}{m}} \leq \min_{\boldsymbol{p} \in \mathcal{M}_m} \max_{i \in \{1, 2, \dots, r\}} \|\boldsymbol{p}(\lambda_i)\|^{\frac{1}{m}}.$$

► Assume that the Hessian H has r distinct eigenvalues in decreasing order: λ₁ > λ₂ > · · · > λ_r

Then

$$\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) \leq \min_{\boldsymbol{p} \in \mathcal{M}_m} \|\boldsymbol{p}(\mathbf{H})\|^{\frac{1}{m}} \leq \min_{\boldsymbol{p} \in \mathcal{M}_m} \max_{i \in \{1, 2, \dots, r\}} \|\boldsymbol{p}(\lambda_i)\|^{\frac{1}{m}}.$$

► Example I: $\rho^{(m)}(\mathbf{H}, \mathbf{g}) \leq (\prod_{i=1}^{m} \lambda_i)^{\frac{1}{m}}$ when m < r, $\rho^{(m)}(\mathbf{H}, \mathbf{g}) = 0$ when $m \geq r$

• If the rank of Hessian is at most m-1, then $ho^{(m)}({f H},{m g})=0$

Assume that the Hessian H has r distinct eigenvalues in decreasing order: λ₁ > λ₂ > · · · > λ_r

Then

$$\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) \leq \min_{\boldsymbol{p} \in \mathcal{M}_m} \|\boldsymbol{p}(\mathbf{H})\|^{\frac{1}{m}} \leq \min_{\boldsymbol{p} \in \mathcal{M}_m} \max_{i \in \{1, 2, \dots, r\}} \|\boldsymbol{p}(\lambda_i)\|^{\frac{1}{m}}.$$

• Example I: $\rho^{(m)}(\mathbf{H}, \mathbf{g}) \leq (\prod_{i=1}^{m} \lambda_i)^{\frac{1}{m}}$ when m < r, $\rho^{(m)}(\mathbf{H}, \mathbf{g}) = 0$ when $m \geq r$ • If the rank of Hessian is at most m - 1, then $\rho^{(m)}(\mathbf{H}, \mathbf{g}) = 0$

Example II: Assume that all the eigenvalues of H lie within [0, Δ] ∪ [L₁ − Δ, L₁] for some L₁ > Δ > 0. Then, when m is even, we have

$$\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) \leq 2^{1/m} \sqrt{\Delta(L_1 - \Delta)}/2$$



Introduction

The Proposed Method

③ Convergence Analysis

4 Numerical Results

Logistic Regression Problems



- The convergence path of Krylov-CRN remains almost unchanged as d increases
- SSCN converges slower as d increases

Logistic Regression Problems



▶ When *d* is large, Krylov CRN converges much faster than the others

Ruichen Jiang, Parameswaran Raman, Shoham Sabach, Aryan Mokhtari, Mingyi Hong, and Volkan Cevher. Krylov Cubic Regularized Newton: A Subspace Second-Order Method with Dimension-Free Convergence Rate, 2024. arXiv: 2401.03058 [math.OC]

Thank you!