# 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
TL;DR		Krylov CRN		Hessian Eigenspectrum	
<ul> <li>We propose a subspace cubic regularized convex minimization</li> <li>The first dimension-independent global subspace methods, where <i>m</i> is the subspace methods, where <i>m</i> is the subspace institute.</li> <li>Key idea: use the Krylov subspace institute.</li> <li>Empirically, our method with <i>m</i> = 10 subspace methods in logistic regression.</li> </ul>	and rate of $O(\frac{1}{mk} + \frac{1}{k^2})$ among bspace dimension tead of a random subspace outperforms existing n problems with $d = 10^6$	<ul> <li>Input: Initial point x<sub>0</sub> ∈ ℝ<sup>d</sup>, subspace dimension m, regularization parameter M &gt; 0</li> <li>for k = 0, 1,, do</li> <li>(V<sub>k</sub>, ğ<sub>k</sub>, H <sub>k</sub>) ← LANCZOS(H<sub>k</sub>, g<sub>k</sub>; m)</li> <li>Solve the cubic subproblem</li> <li>z<sub>k</sub> = argmin z∈ℝ<sup>m</sup> { ğ<sub>k</sub><sup>T</sup> z + 1/2 z<sup>T</sup> H <sub>k</sub> z + M/6   z  <sup>3</sup> },</li> </ul>		► The presented bounds rely on $L_1$ , the Hessian's largest eigenvalue ► 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, \boldsymbol{g}_i)\}, $ where $ \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}} $	
Cubic Regularized Newton MethodThe unconstrained minimization problem $\min_{\boldsymbol{x} \in \mathbb{R}^d} f(\boldsymbol{x})$		<ul> <li>Update x<sub>k+1</sub> = x<sub>k</sub> + V<sub>k</sub>z<sub>k</sub></li> <li>► end for</li> <li>Per-iteration computational cost <ul> <li>Lanczos iteration: m HVPs ⇒ O(md)</li> <li>Solving the cubic subproblem: O(m)</li> </ul> </li> </ul>		<ul> <li>For a monic polynomial p(x) = x<sup>m</sup> - ∑<sub>i=0</sub><sup>m-1</sup> c<sub>i</sub>x<sup>i</sup>, define p(H) = H<sup>m</sup> - ∑<sub>i=0</sub><sup>m-1</sup> c<sub>i</sub>H<sup>i</sup></li> <li>Let M<sub>m</sub> be the set of monic polynomials of degree m. Then ρ<sup>(m)</sup>(H, g) = min p∈M<sub>m</sub>    p(H) g/   <sup>1/m</sup> ≤ min p(H)   <sup>1/m</sup></li> </ul>	

- *f* : ℝ<sup>w</sup> → ℝ is convex *f* is bounded from below and has bounded level-sets • The Hessian of f is Lipchitz, i.e.,  $\|
  abla^2 f(oldsymbol{x}) - 
  abla^2 f(oldsymbol{y})\| \leq L_2 \|oldsymbol{x} - oldsymbol{y}\|$ ,  $orall oldsymbol{x}, oldsymbol{y} \in \mathbb{R}^d$
- ► Let  $\boldsymbol{g}_k := \nabla f(\boldsymbol{x}_k)$  and  $\mathbf{H}_k := \nabla^2 f(\boldsymbol{x}_k)$ Cubic regularized Newton (CRN) [Griewank'81; Nesterov-Polyak'06]  $\boldsymbol{s}_{k} = \operatorname*{argmin}_{\boldsymbol{s} \in \mathbb{R}^{d}} \left\{ \boldsymbol{g}_{k}^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \mathbf{H}_{k} \boldsymbol{s} + \frac{M}{6} \|\boldsymbol{s}\|^{3} \right\}$  $oldsymbol{x}_{k+1} = oldsymbol{x}_k + oldsymbol{s}_k$
- f convex:  $f(\boldsymbol{x}_k) f^* = \mathcal{O}(1/k^2)$
- *f* strongly convex: a superlinear convergence rate

# **Prior Works: Stochastic Subspace CRN**

- ► The drawback of CRN: high memory and computational costs
- Computing & storing the Hessian:  $\mathcal{O}(d^2)$
- Solving a cubic subproblem:  $\mathcal{O}(d^3)$
- Subspace methods: executing 2nd-order updates in a subspace  $\mathcal{V}_k$ of dim  $m \ll d$  [Doikov-Richtárik'18; Gower et al.'19; Hanzely et al.'20]

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

▶ Implementation: let  $\mathbf{V}_k \in \mathbb{R}^{d \times m}$  whose columns form an orthonormal basis for  $\mathcal{V}_k$ . Then

$$\boldsymbol{z}_{k} = \operatorname*{argmin}_{\boldsymbol{z} \in \mathbb{R}^{m}} \left\{ \tilde{\boldsymbol{g}}_{k}^{\top} \boldsymbol{z} + \frac{1}{2} \boldsymbol{z}^{\top} \tilde{\boldsymbol{H}}_{k} \boldsymbol{z} + \frac{M}{6} \|\boldsymbol{z}\|^{3} \right\}$$

 $(\mathbf{V}, \tilde{\boldsymbol{g}}, \tilde{\mathbf{H}}) = \text{Lanczos}(\mathbf{H}, \boldsymbol{g}; m)$ **▶ Input:**  $\mathbf{H} \in \mathbb{R}^{d \times d}$ ,  $\boldsymbol{g} \in \mathbb{R}^{d}$ , and the dimension m► Initialize:  $v_1 = g/||g||$ ,  $\beta_1 = 0$ ,  $v_0 = 0$ ▶ for j = 1, 2, ..., m do // one Hessian-vector product (HVP) •  $\boldsymbol{w}_{j} \leftarrow \mathbf{H} \boldsymbol{v}_{j} - \beta_{j} \boldsymbol{v}_{j-1}$ •  $\alpha_i \leftarrow \boldsymbol{w}_i^\top \boldsymbol{v}_i$ •  $\boldsymbol{w}_i \leftarrow \boldsymbol{w}_i - \alpha_i \boldsymbol{v}_i$ •  $\beta_{j+1} \leftarrow \| \boldsymbol{w}_j \|_2$ •  $\boldsymbol{v}_{j+1} \leftarrow \boldsymbol{w}_j / \beta_{j+1}$  $\bullet \text{ output: } \mathbf{V} = [\boldsymbol{v}_1, \dots, \boldsymbol{v}_m], \ \tilde{\boldsymbol{g}} = \|\boldsymbol{g}\|\boldsymbol{e}_1, \ \tilde{\mathbf{H}} = \begin{bmatrix} \alpha_1 & \beta_2 \\ \beta_2 & \alpha_2 & \beta_3 \\ \beta_3 & \ddots & \ddots \\ & \beta_3 & \ddots & \ddots \\ & & \ddots & \ddots & \beta_{m-1} \\ & & & \beta_{m-1} & \alpha_{m-1} & \beta_m \\ & & & & \beta_m & \alpha_m \end{bmatrix}$ 

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

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

**Example I:**  $\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) \leq \left(\prod_{i=1}^{m} \lambda_i\right)^{\frac{1}{m}}$  when  $m \leq r$ • If the rank of Hessian  $\leq m-1$ , then  $\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) = 0$ 

**Example II**: Assume that all the eigenvalues of **H** lie within  $[0,\Delta] \cup [L_1 - \Delta, L_1]$  for some  $L_1 > \Delta > 0$ . When m is even,  $\rho^{(m)}(\mathbf{H}, \boldsymbol{g}) \leq 2^{1/m} \sqrt{\Delta(L_1 - \Delta)}/2$ 

### **Convergence Analysis**

```
► Key inequality in the analysis of classical CRN:
               f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) + \boldsymbol{g}_k^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \mathbf{H}_k \boldsymbol{s} + \frac{L_2}{6} \|\boldsymbol{s}\|^3, \quad \forall \boldsymbol{s} \in \mathbb{R}^d
► For subspace CRN, we instead have
                 f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) + \boldsymbol{g}_k^{\top} \boldsymbol{s} + \frac{1}{2} \boldsymbol{s}^{\top} \mathbf{H}_k \boldsymbol{s} + \frac{L_2}{6} \|\boldsymbol{s}\|^3, \quad \forall \boldsymbol{s} \in \boldsymbol{\mathcal{V}}_k
```

#### **Logistic Regression Problems** Plot by time: Plot by iteration: w8a (*n* = 49749, *d* = 300) w8a (*n* = 49749, *d* = 300) --- CRN SSCN (m = 10) SSCN (m = Krylov CRN (m = 1

# $\boldsymbol{s}_k = \boldsymbol{V}_k \boldsymbol{z}_k, \quad \boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{s}_k$ where $\tilde{\boldsymbol{g}}_k = \mathbf{V}_k^\top \boldsymbol{g}_k \in \mathbb{R}^m$ and $\tilde{\mathbf{H}}_k = \mathbf{V}_k^\top \mathbf{H}_k \mathbf{V}_k \in \mathbb{R}^{m \times m}$

- $\blacktriangleright$  The existing works choose a random subspace  $\mathcal{V}_k$ 
  - 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:  $\mathcal{O}(m^2)$  in some special cases
- Solving a cubic subproblem:  $\mathcal{O}(m^3)$
- ► ... but a much slower convergence rate

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

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

# **Our Contributions: Krylov CRN**

- We propose the Krylov CRN method, where  $\mathcal{V}_k$  is chosen as the Krylov subspace span $\{g_k, H_kg_k, \ldots, H_k^{m-1}g_k\}$
- ► Can be implemented using the Lanczos method, with 1 gradient evaluation and m Hessian-vector products (HVPs) per iteration In the convex case, we prove a dimension free convergence rate



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

Main Results

**Theorem (Convex Setting)** 

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

- Achieving accuracy  $\epsilon$  requires  $\mathcal{O}\left(\frac{1}{m\epsilon} + \frac{1}{\sqrt{\epsilon}}\right)$  iterations
- ► 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

Theorem (Strongly Convex Setting)







## ► In terms of iterations:

- CRN

→ SSCN (m = 50) SSCN (m = 100) SSCN (m = 500)

• The convergence path of Krylov CRN remains almost unchanged as d increases from 300 to  $10^6$ 

In the convex case, we prove a dimension-free convergence rate					
Methods	Per-iteration cost	Convergence rate			
CRN [Nesterov-Polyak'06]	${\cal O}(d^3)$	$\mathcal{O}(rac{1}{k^2})$			
SSCN [Hanzely et al.'20]	${\cal O}(m^3)^*$	$\mathcal{O}(rac{d-m}{m}\cdotrac{1}{k}+rac{d^2}{m^2}\cdotrac{1}{k^2})$			
Krylov CRN (ours)	$\mathcal{O}(md)^{**}$	$\mathcal{O}(rac{1}{mk}+rac{1}{k^2})$			

\*Assume computing the subspace gradient is  $\mathcal{O}(m)$  and the subspace Hessian is  $\mathcal{O}(m^2)$ \*\*Assume the cost of Hessian-vector product evaluations is  $\mathcal{O}(d)$ 

► When the Hessian spectrum possesses a certain structure, our method achieves a faster rate

•• \_\_\_\_

Let  $\{x_k\}$  be generated by Krylov CRN. Then the number of iterations needed to reach  $\delta_k := f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) \leq \epsilon$  is upper bounded by  $k = \mathcal{O}\left(\left(\frac{L_1}{m\mu} + 1\right)\log\frac{\delta_0}{\epsilon} + \frac{\sqrt{L_2}\delta_0^{0.25}}{\mu^{0.75}}\right)$ 

► In comparison, SSCN requires  $\mathcal{O}\left(\left(\frac{d-mL_1}{m}+\frac{d}{m}\right)\log\frac{1}{\epsilon}\right)$ • Again, when  $m \ll d$ , we shave a factor of d

• SSCN converges slower as d increases

► In terms of time:

• As d increases, the gap between Krylov CRN and other methods becomes wider



Take a picture to download the full paper

