# High Probability Complexity Bounds for Trust-Region Methods with Noisy Oracles

Liyuan Cao 曹立元

Beijing International Center for Mathematical Research, Peking University
北京大学北京国际数学研究中心

April 8, 2023

Albert S. Berahas
University of Michigan



Katya Scheinberg
Cornell University

Liyuan Cao, Albert S. Berahas, and Katya Scheinberg. First-and second-order high probability complexity bounds for trust-region methods with noisy oracles. *arXiv preprint arXiv:2205.03667*, 2022.

$$\min_{x \in \mathbb{R}^n} \phi(x)$$

$\phi$ follows common assumptions

$$\phi(x) \geq \hat{\phi} \text{ for all } x \in \mathbb{R}^n,$$
$$\|\nabla\phi(x) - \nabla\phi(y)\| \leq L_1\|x - y\| \text{ for all } (x, y) \in \mathbb{R}^n \times \mathbb{R}^n,$$

but we only have access to

$$\left\{ \begin{array}{l} f_k \\ g_k \\ H_k \end{array} \right. \quad \text{instead of} \quad \left\{ \begin{array}{l} \phi(x_k) \\ \nabla\phi(x_k) \\ \nabla^2\phi(x_k). \end{array} \right.$$

$$\min_{x \in \mathbb{R}^n} \phi(x)$$

$\phi$ follows common assumptions

$$\phi(x) \geq \hat{\phi} \text{ for all } x \in \mathbb{R}^n,$$
$$\|\nabla\phi(x) - \nabla\phi(y)\| \leq L_1 \|x - y\| \text{ for all } (x, y) \in \mathbb{R}^n \times \mathbb{R}^n,$$

but we only have access to

$$\left\{ \begin{array}{l} f_k \\ g_k \\ H_k \end{array} \right. \quad \text{instead of} \quad \left\{ \begin{array}{l} \phi(x_k) \\ \nabla\phi(x_k) \\ \nabla^2\phi(x_k). \end{array} \right.$$

A line of work:

algorithm TR method modified to handle noise.

noise Weaker assumptions in more recent work.

result Stronger results in more recent work.

---

**Algorithm: Modified First-Order Trust-Region Method**

---

**Inputs:** Starting point $x_0$, initial trust region radius $\delta_0$, tolerance parameter $r$, and hyperparameters $\eta_1 > 0, \eta_2 > 0, \gamma \in (0, 1)$ for controlling the trust region radius.

**for** $k = 0, 1, 2, \ldots$ **do**

1    Build a quadratic model $m_k(x_k + s) = \phi(x_k) + \langle g_k, s \rangle + 0.5 \langle H_k s, s \rangle$

2    Compute $s_k$ by approximately minimizing $m_k$ in $B(x_k, \delta_k)$ so that it satisfies the *Cauchy decrease condition*

$$m_k(x_k) - m_k(x_k + s_k) \geq \frac{1}{2} \|g_k\| \min \left\{ \frac{\|g_k\|}{\|H_k\|}, \delta_k \right\}.$$
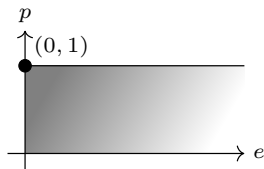
3    Compute

$$\rho_k = \frac{f_k - f_k^+ + r}{m_k(x_k) - m_k(x_k + s_k)}$$
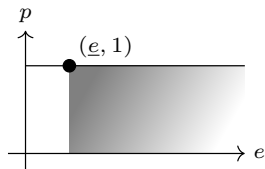
and update $x$ and $\delta$

$$(x_{k+1}, \delta_{k+1}) = \left\{ \begin{array}{ll} (x_k + s_k, \gamma^{-1}\delta_k) & \text{if } \rho_k \geq \eta_1 \text{ and } \|g_k\| \geq \eta_2 \delta_k \\ (x_k + s_k, \gamma\delta_k) & \text{if } \rho_k \geq \eta_1 \text{ and } \|g_k\| < \eta_2 \delta_k \\ (x_k, \gamma\delta_k) & \text{if } \rho_k < \eta_1. \end{array} \right.$$

Let $\varphi^{(j)}\left(x_k, \xi_k^{(j)}, \mathcal{S}_k^{(j)}\right)$ be the $j$th-order oracle that returns an estimate of $\nabla^j \phi(x_k)$ such that for all $(e, p) \in \mathcal{S}_k^{(j)} \subseteq [0, \infty) \times [0, 1]$,
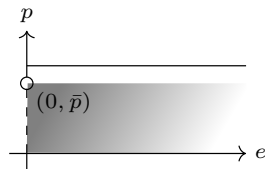
$$\mathbb{P}_{\xi_k^{(j)}} \left\{ \|\varphi^{(j)}\left(x_k, \xi_k^{(j)}, \mathcal{S}^{(j)}\right) - \nabla^j \phi(x_k)\| \le e \Big| \mathcal{F}_k \right\} \ge p$$



**(a)** $\mathcal{S}_k^{(j)} \ni (0, 1)$
exact

**(b)** $\mathcal{S}_k^{(j)} = [\underline{e}, +\infty) \times [0, 1]$
bounded

**(c)** $\mathcal{S}_k^{(j)} = (0, +\infty) \times [0, \bar{p}]$
probabilistically sufficiently accurate

# The Goal of the Analysis

ⓒ convergence
$$\liminf_{k \to \infty} \|\nabla \phi(x_k)\| \le \epsilon$$

ⓔ expected complexity
$$\mathbb{E} \min\{k : \|\nabla \phi(X_k)\| \le \epsilon\} = \mathcal{O}(1/\epsilon^2)$$

ⓗ high probability convergence

$\mathbb{P}\{\min\{\|\nabla \phi(X_k)\| : \ 0 \le k \le T-1\} < \epsilon\}$
$$\ge \text{a function of } T \text{ the converges to 1 as } T \text{ increase}$$

for some sufficiently large $\epsilon$.

$\mathcal{S}^{(0)} = [0, \infty) \times [0, 1]$ and $\mathcal{S}^{(1)} = (0, \infty) \times [0, \bar{p}_1]$ for sufficiently large $\bar{p}_1$:

- Ⓒ Afonso S Bandeira, Katya Scheinberg, and Luis Nunes Vicente. Convergence of trust-region methods based on probabilistic models. *SIAM Journal on Optimization*, 24(3):1238–1264, 2014.

- Ⓗ Serge Gratton, Clement W Royer, Luis N Vicente, and Zaikun Zhang. Complexity and global rates of trust-region methods base on probabilistic models. *IMA Journal of Numerical Analysis*, 38(3):1579-1597, 2018.

$\mathcal{S}^{(j)} = (0, \infty) \times [0, \bar{p}_j]$ for sufficiently large $\bar{p}_j$, $j = 0, 1$:

- Ⓒ Ruobing Chen, Matt Menickelly, and Katya Scheinberg. Stochastic optimization using a trust-region method and random models. *Mathematical Programming*, 169(2):447–487, 2018.

$\mathcal{S}^{(j)} = (0, \infty) \times [0, \bar{p}_j]$ for sufficiently large $\bar{p}_j$, $j = 0, 1, 2$ and $\mathbb{E}_{\xi_0}|f_k - \phi(x_k)| \leq C_0$:

- Ⓔ Jose Blanchet, Coralia Cartis, Matt Menickelly, and Katya Scheinberg. Convergence rate analysis of a stochastic trust-region method via supermartingales. *INFORMS Journal on Optimization*, 1(2):92–119, 2019.

$\mathcal{S}^{(0)} = [\epsilon_f, \infty) \times [0, 1]$ and $\mathcal{S}^{(1)} = [\epsilon_g, \infty) \times [0, 1]$:

- Ⓔ Shigeng Sun and Jorge Nocedal. A trust region method for the optimization of noisy functions. *arXiv preprint arXiv:2201.00973*, 2022.
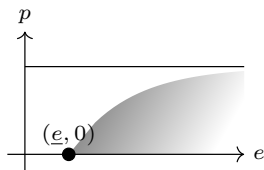
We assume for all $k$:

$$\left.\begin{array}{l} \mathbb{P}\left\{|f_k - \phi(x_k)| \leq e\right\} \\ \mathbb{P}\left\{|f_k^+ - \phi(x_k + s_k)| \leq e\right\} \end{array}\right\} \geq \exp(a(\epsilon_f - e)) \qquad \text{unbounded noise,}$$

$$\mathbb{P}\left\{\|g_k - \nabla\phi(x_k)\| \leq \kappa_{\mathrm{eg}}\delta_k + \epsilon_g\right\} \geq p_1 \qquad \text{irreducible noise,}$$

$\|H_k\| \leq \kappa_{\mathrm{bhm}}$ for some constant $\kappa_{\mathrm{bhm}}$ (bound on hessian of model).



**(a)** $\mathcal{S}_k^{(0)} = \{(e,p) : e \geq \underline{e} = \epsilon_f, p \leq 1 - \exp(a(\underline{e} - e))\}$

**(b)** $\mathcal{S}_k^{(1)} = (\underline{e}, +\infty) \times [0, \bar{p}]$ with $\underline{e} = \kappa_{\mathrm{eg}}\delta_k + \epsilon_g$ and $\bar{p} = p_1$

## Algorithm: Modified First-Order Trust-Region Method

**Inputs:** Starting point $x_0$, initial trust region radius $\delta_0$, tolerance parameter $r$, and hyperparameters $\eta_1 > 0, \eta_2 > 0$, $\gamma \in (0,1)$ for controlling the trust region radius.

**for** $k = 0, 1, 2, \ldots$ **do**

1.     Build a quadratic model $m_k(x_k + s) = \phi(x_k) + \langle g_k, s \rangle + 0.5 \langle H_k s, s \rangle$

2.     Compute $s_k$ by approximately minimizing $m_k$ in $B(x_k, \delta_k)$ so that it satisfies the *Cauchy decrease condition*

$$m_k(x_k) - m_k(x_k + s_k) \geq \frac{1}{2} \|g_k\| \min \left\{ \frac{\|g_k\|}{\|H_k\|}, \delta_k \right\}.$$

3.     Compute

$$\rho_k = \frac{f_k - f_k^+ + r}{m_k(x_k) - m_k(x_k + s_k)}$$

and update $x$ and $\delta$

$$(x_{k+1}, \delta_{k+1}) = \left\{ \begin{array}{ll} (x_k + s_k, \gamma^{-1}\delta_k) & \text{if } \rho_k \geq \eta_1 \text{ and } \|g_k\| \geq \eta_2\delta_k \\ (x_k + s_k, \gamma\delta_k) & \text{if } \rho_k \geq \eta_1 \text{ and } \|g_k\| < \eta_2\delta_k \\ (x_k, \gamma\delta_k) & \text{if } \rho_k < \eta_1. \end{array} \right.$$

| random variables: | $X_k$ | $X_k^+$ | $\mathcal{E}_k$ | $\mathcal{E}_k^+$ | |
|---|---|---|---|---|---|
| realizations: | $x_k$ | $x_k + s_k$ | $\|f_k - \phi(x_k)\|$ | $\|f_k^+ - \phi(x_k + s_k)\|$ | |
| random variables: | $M_k$ | $\nabla M_k$ | $\nabla^2 M_k$ | $\Delta_k$ | $\rho_k$ |
| realizations: | $m_k$ | $g_k$ | $H_k$ | $\delta_k$ | $\rho_k$ |

Define

$$I_k = \mathbb{1}\{\|\nabla M_k - \nabla\phi(X_k)\| \leq \kappa_{\text{eg}}\Delta_k + \epsilon_g\} \qquad \text{gradient sufficiently accurate}$$

$$J_k = \mathbb{1}\{\mathcal{E}_k + \mathcal{E}_k^+ \leq r\} \qquad \text{zeroth-order noise compensated}$$

$$\Lambda_k = \mathbb{1}\{\Delta_k > \bar{\Delta}\} \qquad \text{large TR radius}$$

$$\Theta_k = \mathbb{1}\{\rho_k \geq \eta_1 \text{ and } \|\nabla M_k\| \geq \eta_2\Delta_k\} \qquad \text{successful step}$$

$$\Theta_k' = \mathbb{1}\{\rho_k \geq \eta_1\} \qquad \text{accepted step}$$

where $\bar{\Delta} = C_1 \min_{0 \leq k \leq T-1} \|\nabla\phi(X_k)\| - C_2\epsilon_g$.

# Classification of Iterations

| | $I_k=1, J_k=1$ | | | $I_k=1, J_k=0$ | | | $I_k=0, J_k=1$ | | | $I_k=0, J_k=0$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ★ | ✓ | ✗ | ★ | ✓ | ✗ | ★ | ✓ | ✗ | ★ | ✓ | ✗ |
| $\Delta_k \in (\bar{\Delta}, \infty)$ | 1 | 4 | 5 | 6 | 9 | 11 | 13 | 16 | 18 | 20 | 23 | 25 |
| $\Delta_k \in (\gamma\bar{\Delta}, \bar{\Delta}]$ | 2 | | | 7 | | | 14 | | | 21 | | |
| $\Delta_k \in (0, \gamma\bar{\Delta}]$ | 3 | | | 8 | 10 | 12 | 15 | 17 | 19 | 22 | 24 | 26 |

## Lemma (sufficient condition for successful step)

*If $I_k J_k = 1$ and $\Lambda_k = 0$ then $\Theta_k = 1$.*

## Lemma (**progress made in each iteration**)

*Let $h(\delta) = C_3 \delta^2$. Then we have*

$$\phi(X_k) - \phi(X_{k+1}) \geq \begin{cases} h(\Delta_k) - \mathcal{E}_k - \mathcal{E}_k^+ - r & \text{if } \Theta_k = 1 \text{ (successful)} \\ -\mathcal{E}_k - \mathcal{E}_k^+ - r & \text{if } \Theta_k' = 1 \text{ (accepted)} \\ 0 & \text{if } \Theta_k' = 0 \text{ (rejected)}. \end{cases}$$

## Lemma (**total progress**)

$$h(\gamma \bar{\Delta}) \sum_{k=0}^{T-1} \Theta_k \Lambda_k \leq \phi(x_0) - \hat{\phi} + \sum_{k=0}^{T-1} \Theta_k' \left( \mathcal{E}_k + \mathcal{E}_k^+ + r \right).$$

# Lemma of Total Progress

### Lemma (**total loss**)

*For any $t \geq 0$,*

$$\mathbb{P}\left\{\sum_{k=0}^{T-1}\left(\mathcal{E}_k + \mathcal{E}_k^+ + r\right) \geq T(4/a + 2\epsilon_f + r) + t\right\} \leq \exp\left(-\frac{a}{4}t\right).$$

Let $t = rT$.

### Lemma (**total progress**)

$$h(\gamma\bar{\Delta})\sum_{k=0}^{T-1}\Theta_k\Lambda_k \leq \phi(x_0) - \hat{\phi} + \sum_{k=0}^{T-1}\Theta_k'\left(\mathcal{E}_k + \mathcal{E}_k^+ + r\right)$$

$$< \phi(x_0) - \hat{\phi} + T(4/a + 2\epsilon_f + 2r)$$

*with probability at least $1 - \exp\left(-\frac{ar}{4}T\right)$.*

$$h(\gamma\bar{\Delta}) = \gamma^2 C_3\left(C_1 \min_{0 \leq k \leq T-1}\|\nabla\phi(X_k)\| - C_2\epsilon_g\right)^2$$
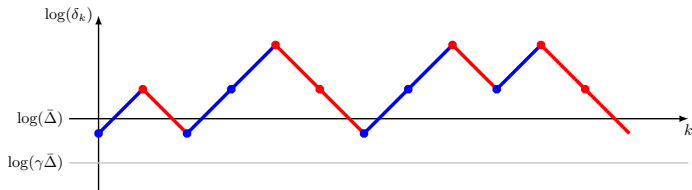
| | $I_k=1, J_k=1$ | | | $I_k=1, J_k=0$ | | | $I_k=0, J_k=1$ | | | $I_k=0, J_k=0$ | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | ★ | ✓ | ✗ | ★ | ✓ | ✗ | ★ | ✓ | ✗ | ★ | ✓ | ✗ |
| $\Delta_k \in (\bar{\Delta}, \infty)$ | 1 | 4 | 5 | 6 | 9 | 11 | 13 | 16 | 18 | 20 | 23 | 25 |
| $\Delta_k \in (\gamma\bar{\Delta}, \bar{\Delta}]$ | 2 | | | 7 | 10 | 12 | 14 | 17 | 19 | 21 | 24 | 26 |
| $\Delta_k \in (0, \gamma\bar{\Delta}]$ | 3 | | | 8 | | | 15 | | | 22 | | |

**Lemma (total progress)**

$$h(\gamma\bar{\Delta}) \sum_{k=0}^{T-1} \Theta_k \Lambda_k \leq \phi(x_0) - \hat{\phi} + T(4/a + 2\epsilon_f + 2r).$$
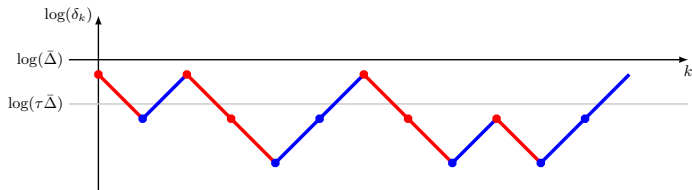
# Ups and Downs of the Radius



$$\sum_{k=0}^{T-1}(1-\Theta_k)\Lambda_k < \sum_{k=0}^{T-1}\Theta_k\Lambda_k + \min\left\{\log_\gamma\left(\frac{\delta_0}{\bar{\Delta}}\right), 0\right\} + 1$$

# Downs and Ups of the Radius



$$\sum_{k=0}^{T-1} \Theta_k (1 - \Lambda_k) < \sum_{k=0}^{T-1} (1 - \Theta_k)(1 - \Lambda_k) + \min\left\{\log_\gamma\left(\frac{\bar{\Delta}}{\delta_0}\right), 0\right\} + 1$$

# Iterations with Sufficiently Accurate Gradient Estimate

| | $I_k = 1, J_k = 1$ | | | $I_k = 1, J_k = 0$ | | | $I_k = 0, J_k = 1$ | | | $I_k = 0, J_k = 0$ | | |
| | ★ | ✓ | ✗ | ★ | ✓ | ✗ | ★ | ✓ | ✗ | ★ | ✓ | ✗ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\Delta_k \in (\bar{\Delta}, \infty)$ | 1 | 4 | 5 | 6 | 9 | 11 | 13 | 16 | 18 | 20 | 23 | 25 |
| $\Delta_k \in (\gamma\bar{\Delta}, \bar{\Delta}]$ | 2 | | | 7 | | | 14 | 17 | 19 | 21 | 24 | 26 |
| $\Delta_k \in (0, \gamma\bar{\Delta}]$ | 3 | | | 8 | 10 | 12 | 15 | | | 22 | | |

## Lemma

*Assume $\mathbb{P}\{I_k = 1 \mid \mathcal{F}_k\} \geq p_1$ holds. By Azuma-Hoeffding inequality, for any positive integer $T$ and any $\hat{p}_1 \in [0, p_1]$ we have*

$$\mathbb{P}\left\{ \sum_{k=0}^{T-1} I_k > \hat{p}_1 T \right\} \geq 1 - \exp\left( -\frac{(1 - \hat{p}_1/p_1)^2}{2} T \right).$$

# Iterations with Sufficiently Accurate Function Evaluation



|  | $I_k=1, J_k=1$ | | | $I_k=1, J_k=0$ | | | $I_k=0, J_k=1$ | | | $I_k=0, J_k=0$ | | |
|---|★|✓|✗|★|✓|✗|★|✓|✗|★|✓|✗|
| $\Delta_k \in (\bar{\Delta}, \infty)$ | 1 | 4 | 5 | 6 | 9 | 11 | 13 | 16 | 18 | 20 | 23 | 25 |
| $\Delta_k \in (\gamma\bar{\Delta}, \bar{\Delta}]$ | 2 | | | 7 | | | 14 | | | 21 | | |
| $\Delta_k \in (0, \gamma\bar{\Delta}]$ | 3 | | | 8 | 10 | 12 | 15 | 17 | 19 | 22 | 24 | 26 |

## Lemma

*Assume both $\mathbb{P}\{\mathcal{E}_k > t\}$ and $\mathbb{P}\{\mathcal{E}_k^+ > t\}$ are $\leq \exp(a(\epsilon_f - t))$. Let $p_0 = 1 - 2\exp(a[\epsilon_f - r/2])$. For any positive integer $T$ and any $\hat{p}_0 \in [0, p_0]$, we have*

$$\mathbb{P}\left\{ \sum_{k=0}^{T-1} J_k > \hat{p}_0 T \right\} \geq 1 - \exp\left( -\frac{(1 - \hat{p}_0/p_0)^2}{2} T \right).$$

$$\sum_{k=0}^{T-1} (1-\Theta_k)\Lambda_k < \sum_{k=0}^{T-1} \Theta_k\Lambda_k + \min\left\{\log_\gamma\left(\frac{\delta_0}{\bar{\Delta}}\right), 0\right\} + 1$$

$$\sum_{k=0}^{T-1} \Theta_k(1-\Lambda_k) < \sum_{k=0}^{T-1} (1-\Theta_k)(1-\Lambda_k) + \min\left\{\log_\gamma\left(\frac{\bar{\Delta}}{\delta_0}\right), 0\right\} + 1$$

$$\mathbb{P}\left\{\sum_{k=0}^{T-1} I_k > \hat{p}_1 T\right\} \geq 1 - \exp\left(-\frac{(1-\hat{p}_1/p_1)^2}{2}T\right)$$

$$\mathbb{P}\left\{\sum_{k=0}^{T-1} J_k > \hat{p}_0 T\right\} \geq 1 - \exp\left(-\frac{(1-\hat{p}_0/p_0)^2}{2}T\right)$$

$$\Downarrow$$

$$\mathbb{P}\left\{\sum_{k=0}^{T-1} \Theta_k\Lambda_k > \left(\hat{p}_0 + \hat{p}_1 - \frac{3}{2}\right)T - \frac{1}{2}\left|\log_\gamma\frac{\bar{\Delta}}{\delta_0}\right| - \frac{1}{2}\right\}$$

$$\geq 1 - \exp\left(-\frac{(1-\hat{p}_1/p_1)^2}{2}T\right) - \exp\left(-\frac{(1-\hat{p}_0/p_0)^2}{2}T\right)$$

# Main Result

Let $t = rT$.

## Theorem

*Let assumptions hold. Given any $\epsilon > \sqrt{\frac{4\epsilon_f + 8/a + 2r}{C_3\gamma^2 C_1^2(2p_0 + 2p_1 - 3)}} + \frac{C_2}{C_1}\epsilon_g$, we have*

$$\mathbb{P}\left\{\min\{\|\nabla\phi(X_k)\| : \ 0 \le k \le T - 1\} \le \epsilon\right\} \ge$$
$$1 - \exp\left(-\frac{(1 - \hat{p}_1/p_1)^2}{2}T\right) - \exp\left(-\frac{(1 - \hat{p}_0/p_0)^2}{2}T\right) - \exp\left(-\frac{ar}{4}T\right)$$

*for any $\hat{p}_0$ and $\hat{p}_1$ such that $\hat{p}_0 + \hat{p}_1 \in \left(\frac{3}{2} + \frac{2\epsilon_f + 4/a + r}{C_3\gamma^2(C_1\epsilon - C_2\epsilon_g)^2}, p_0 + p_1\right]$, any $t \ge 0$, and any*

$$T \ge \left(\hat{p}_0 + \hat{p}_1 - \frac{3}{2} - \frac{2\epsilon_f + 4/a + 2r}{C_3\gamma^2(C_1\epsilon - C_2\epsilon_g)^2}\right)^{-1}$$
$$\left[\frac{\phi(x_0) - \hat{\phi}}{C_3\gamma^2(C_1\epsilon - C_2\epsilon_g)^2} + \frac{1}{2}\left|\log_\gamma \frac{C_1\epsilon - C_2\epsilon_g}{\delta_0}\right| + \frac{1}{2}\right] = \bar{\mathcal{O}}(\epsilon^{-2}).$$

# Other Results

- Analyses under bounded noise assumption.
- Second-order TR method and analysis.
- Numerically testing the strength of the theoretical results.
- Experimenting with different values for $r$.