

Stochastic Second-Order Optimization Methods

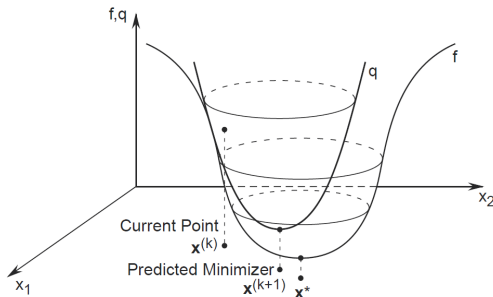
Part I: Convex

Fred Roosta

School of Mathematics and Physics
University of Queensland

Iterative Scheme

$$\mathbf{y}^{(k)} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ (\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{g}^{(k)} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{H}^{(k)} (\mathbf{x} - \mathbf{x}^{(k)}) \right\}$$
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k (\mathbf{y}^{(k)} - \mathbf{x}^{(k)})$$



Iterative Scheme

$$\mathbf{y}^{(k)} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ (\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{g}^{(k)} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{H}^{(k)} (\mathbf{x} - \mathbf{x}^{(k)}) \right\}, \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k (\mathbf{y}^{(k)} - \mathbf{x}^{(k)})$$

- **Newton:** $\mathbf{g}^{(k)} = \nabla F(\mathbf{x}^{(k)})$ & $\mathbf{H}^{(k)} = \nabla^2 F(\mathbf{x}^{(k)})$
- **Gradient Descent:** $\mathbf{g}^{(k)} = \nabla F(\mathbf{x}^{(k)})$ & $\mathbf{H}^{(k)} = \mathbf{I}$
- **Frank-Wolfe:** $\mathbf{g}^{(k)} = \nabla F(\mathbf{x}^{(k)})$ & $\mathbf{H}^{(k)} = \mathbf{0}$
- **(mini-batch) SGD:** $\mathcal{S}_g \subset \{1, 2, \dots, n\} \implies \mathbf{g}^{(k)} = \frac{1}{|\mathcal{S}_g|} \sum_{j \in \mathcal{S}_g} \nabla f_j(\mathbf{x}^{(k)})$ & $\mathbf{H}^{(k)} = \mathbf{I}$
- **Sub-Sampled Newton:**

Hessian Sub-Sampling

$$\mathbf{g}^{(k)} = \nabla F(\mathbf{x}^{(k)})$$

$$\mathcal{S}_g \subset \{1, 2, \dots, n\} \implies \mathbf{H}^{(k)} = \frac{1}{|\mathcal{S}_H|} \sum_{j \in \mathcal{S}_H} \nabla^2 f_j(\mathbf{x}^{(k)})$$

Gradient and Hessian Sub-Sampling

$$\mathcal{S}_g \subset \{1, 2, \dots, n\} \implies \mathbf{g}^{(k)} = \frac{1}{|\mathcal{S}_g|} \sum_{j \in \mathcal{S}_g} \nabla f_j(\mathbf{x}^{(k)})$$

$$\mathcal{S}_H \subset \{1, 2, \dots, n\} \implies \mathbf{H}^{(k)} = \frac{1}{|\mathcal{S}_H|} \sum_{j \in \mathcal{S}_H} \nabla^2 f_j(\mathbf{x}^{(k)})$$

Sub-sampled Newton's Method

Let $\mathcal{X} = \mathcal{R}^d$, i.e., unconstrained optimization

Iterative Scheme

$$\mathbf{p}^{(k)} \approx \underset{\mathbf{p} \in \mathcal{R}^d}{\operatorname{argmin}} \left\{ \mathbf{p}^T \mathbf{g}^{(k)} + \frac{1}{2} \mathbf{p}^T \mathbf{H}^{(k)} \mathbf{p} \right\},$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha^{(k)} \mathbf{p}^{(k)}$$

Hessian Sub-Sampling

$$\mathbf{g}^{(k)} = \nabla F(\mathbf{x}^{(k)})$$

$$\mathbf{H}^{(k)} = \frac{1}{|\mathcal{S}_H|} \sum_{j \in \mathcal{S}_H} \nabla^2 f_j(\mathbf{x}^{(k)})$$

Sub-sampled Newton's Method

Algorithm Newton's Method with Hessian Sub-Sampling

- 1: **Input:** $\mathbf{x}^{(0)}$
 - 2: **for** $k = 0, 1, 2, \dots$ until termination **do**
 - 3: - $\mathbf{g}^{(k)} = \nabla F(\mathbf{x}^{(k)})$
 - 4: - $\mathcal{S}_H^{(k)} \subseteq \{1, 2, \dots, n\}$
 - 5: - $\mathbf{H}^{(k)} = \frac{1}{|\mathcal{S}_H^{(k)}|} \sum_{j \in \mathcal{S}_H^{(k)}} \nabla^2 f_j(\mathbf{x}^{(k)})$
 - 6: - $\mathbf{H}^{(k)} \mathbf{p}^{(k)} \approx -\mathbf{g}^{(k)}$
 - 7: - Find $\alpha^{(k)}$ that passes Armijo linesearch
 - 8: - Update $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha^{(k)} \mathbf{p}^{(k)}$
 - 9: **end for**
-

Sub-Sampling Hessian

$$\mathbf{H} = \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \nabla^2 f_j(\mathbf{x})$$

Lemma ([Roosta and Mahoney, 2016a])

Suppose $\nabla^2 F(\mathbf{x}) \succeq \gamma \mathbf{I}$ and $0 \preceq \nabla^2 f_i(\mathbf{x}) \preceq L \mathbf{I}$. Given any $0 < \epsilon < 1$, $0 < \delta < 1$, if Hessian is uniformly sub-sampled with

$$|\mathcal{S}| \geq \frac{2\kappa \log(d/\delta)}{\epsilon^2},$$

then

$$\Pr\left(\mathbf{H} \succeq (1 - \epsilon)\gamma \mathbf{I}\right) \geq 1 - \delta.$$

where $\kappa = L/\gamma$.

Sub-sampled Newton's Method

Local convergence, i.e., in a neighborhood of \mathbf{x}^* , and with $\alpha^{(k)} = 1$

Linear-Quadratic Error Recursion

$$\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \leq \xi_1 \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| + \xi_2 \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|^2$$

Sub-sampled Newton's Method

[Erdogdu and Montanari, 2015]

$$[\mathbf{U}_{r+1}, \mathbf{\Lambda}_{r+1}] = \text{TSVD}(\mathbf{H}, r+1)$$

$$\hat{\mathbf{H}}^{-1} = \lambda_{r+1}^{-1} \mathbf{I} + \mathbf{U}_r \left(\mathbf{\Lambda}_r^{-1} - \frac{1}{\lambda_{r+1}} \mathbf{I} \right) \mathbf{U}_r^T$$

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmin}} \left\{ \mathbf{x} - \left(\mathbf{x}^{(k)} - \alpha^{(k)} \hat{\mathbf{H}}^{-1} \nabla F(\mathbf{x}^{(k)}) \right) \right\}$$

Theorem ([Erdogdu and Montanari, 2015])

$$\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \leq \xi_1^{(k)} \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| + \xi_2^{(k)} \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|^2,$$

$$\xi_1^{(k)} = 1 - \frac{\lambda_{\min}(\mathbf{H}^{(k)})}{\lambda_{r+1}(\mathbf{H}^{(k)})} + \frac{L_g}{\lambda_{r+1}(\mathbf{H}^{(k)})} \sqrt{\frac{\log d}{|\mathcal{S}_{\mathbf{H}}^{(k)}|}}, \quad \text{and} \quad \xi_2^{(k)} = \frac{L_{\mathbf{H}}}{2\lambda_{r+1}(\mathbf{H}^{(k)})}.$$

Note: For constrained optimization, the method is based on **two metric gradient projection** \implies starting from an arbitrary point, the algorithm might not recognize, i.e., fail to stop at, an stationary point.

Sub-sampling Hessian

$$\mathbf{H} = \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \nabla^2 f_j(\mathbf{x})$$

Lemma ([Roosta and Mahoney, 2016b])

Given any $0 < \epsilon < 1$, $0 < \delta < 1$, if Hessian is uniformly sub-sampled with

$$|\mathcal{S}| \geq \frac{2\kappa^2 \log(d/\delta)}{\epsilon^2},$$

then

$$\Pr\left((1 - \epsilon) \nabla^2 F(\mathbf{x}) \preceq \mathbf{H}(\mathbf{x}) \preceq (1 + \epsilon) \nabla^2 F(\mathbf{x}) \right) \geq 1 - \delta.$$

Sub-sampled Newton's Method [Roosta and Mahoney, 2016b]

Theorem ([Roosta and Mahoney, 2016b])

With high-probability, we get

$$\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \leq \xi_1 \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| + \xi_2 \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|^2,$$

where

$$\xi_1 = \frac{\epsilon}{(1-\epsilon)} + \left(\sqrt{\frac{\kappa}{1-\epsilon}} \right) \theta, \quad \text{and} \quad \xi_2 = \frac{L_H}{2(1-\epsilon)\gamma}.$$

If $\theta = \epsilon/\sqrt{\kappa}$, then ξ_1 is **problem-independent!** \Rightarrow Can be made **arbitrarily small!**

Sub-sampled Newton's Method

Theorem ([Roosta and Mahoney, 2016b])

Consider any $0 < \rho_0 < \rho < 1$ and $\epsilon \leq \rho_0/(1 + \rho_0)$. If $\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq (\rho - \rho_0)/\xi_2$, and

$$\theta \leq \rho_0 \sqrt{\frac{(1 - \epsilon)}{\kappa}},$$

we get locally **Q-linear** convergence

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \rho \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|, \quad k = 1, \dots, k_0$$

with probability $(1 - \delta)^{k_0}$.

- **Problem-independent** local convergence rate
- By increasing Hessian accuracy, **super-linear** rate is possible

Putting it all together

Theorem ([Roosta and Mahoney, 2016b])

Under certain assumptions, starting at any $\mathbf{x}^{(0)}$, we have

- *linear convergence*
- *after certain number of iterations, we get “**problem-independent**” linear convergence*
- *after certain number of iterations, the step size of $\alpha^{(k)} = 1$ passes Armijo rule for **all** subsequent iterations*

*“Any optimization algorithm for which the **unit step length** works has some wisdom. It is too much of a fluke if the unite step length [accidentally] works.”*

Prof. Jorge Nocedal

IPAM Summer School, 2012

Sub-sampled Newton's Method

Theorem ([Bollapragada et al., 2016])

Suppose

$$\left\| \mathbb{E}_i \left(\nabla^2 f_i(\mathbf{x}) - \nabla^2 F(\mathbf{x}) \right)^2 \right\| \leq \sigma^2.$$

Then,

$$\mathbb{E} \left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \leq \xi_1 \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| + \xi_2 \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|^2,$$

where

$$\xi_1 = \frac{\sigma}{\gamma \sqrt{|\mathcal{S}_H^{(k)}|}} + \kappa \theta, \quad \text{and} \quad \xi_2 = \frac{L_H}{2\gamma}.$$

Exploiting the structure....

Example:

$$F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{a}_i^T \mathbf{x}) \implies \nabla^2 F(\mathbf{x}) = \mathbf{A}^T \mathbf{D} \mathbf{A},$$

where

$$\mathbf{A} \triangleq \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_n^T \end{pmatrix} \in \mathbb{R}^{n \times d}, \quad \mathbf{D}_x \triangleq \frac{1}{n} \begin{pmatrix} \ell''(\mathbf{a}_1^T \mathbf{x}) & & \\ & \ddots & \\ & & \ell''(\mathbf{a}_n^T \mathbf{x}) \end{pmatrix} \in \mathbb{R}^{n \times n}.$$

Sketching [Pilanci and Wainwright, 2017]

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ (\mathbf{x} - \mathbf{x}^{(k)})^T \nabla F(\mathbf{x}^{(k)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(k)})^T \nabla^2 F(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) \right\}$$

$$\nabla^2 F(\mathbf{x}^{(k)}) = \mathbf{B}^{(k)T} \mathbf{B}^{(k)}, \quad \mathbf{B}^{(k)} \in \mathcal{R}^{n \times d}$$

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ (\mathbf{x} - \mathbf{x}^{(k)})^T \nabla F(\mathbf{x}^{(k)}) + \frac{1}{2} \underbrace{\|\mathbf{B}^{(k)}\|}_{n \times d} (\mathbf{x} - \mathbf{x}^{(k)}) \right\}^2$$

$$\nabla^2 F(\mathbf{x}^{(k)}) \approx \mathbf{B}^{(k)T} \mathbf{S}^{(k)T} \mathbf{S}^{(k)} \mathbf{B}^{(k)}, \quad \mathbf{S}^{(k)} \in \mathcal{R}^{s \times n}, \quad d \leq s \ll n$$

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ (\mathbf{x} - \mathbf{x}^{(k)})^T \nabla F(\mathbf{x}^{(k)}) + \frac{1}{2} \underbrace{\|\mathbf{S}^{(k)} \mathbf{B}^{(k)}\|}_{s \times d} (\mathbf{x} - \mathbf{x}^{(k)}) \right\}^2$$

Sketching [Pilanci and Wainwright, 2017]

- Sub-Gaussian sketches
 - well-known concentration properties
 - involve dense/unstructured matrix operations
- Randomized orthonormal systems, e.g., Hadamard or Fourier
 - sub-optimal sample sizes
 - fast matrix multiplication
- Random row sampling
 - Uniform
 - Non-uniform (more on this later)
- Sparse JL sketches

Non-uniform Sampling [Xu et al., 2016]

- **Non-uniformity** among $\nabla^2 f_i \implies \mathcal{O}(n)$ uniform samples!!!
- Find $|\mathcal{S}_H|$ **independent** of n
- **Immune** non-uniformity
- Non-Uniform sampling schemes base on
 - 1 Row norms
 - 2 Leverage scores

LISSA [Agarwal et al., 2017]

- Key idea: use Taylor expansion (Neumann series) to construct an estimator of the Hessian inverse
- $\|\mathbf{A}\| \leq 1, \mathbf{A} \succ 0 \implies \mathbf{A}^{-1} = \sum_{i=0}^{\infty} (\mathbf{I} - \mathbf{A})^i$
- $\mathbf{A}_j^{-1} = \sum_{i=0}^j (\mathbf{I} - \mathbf{A})^i = \mathbf{I} + (\mathbf{I} - \mathbf{A}) \mathbf{A}_{j-1}^{-1}$
- $\lim_{j \rightarrow \infty} \mathbf{A}_j^{-1} = \mathbf{A}^{-1}$
- Uniformly sub-sampled Hessian + Truncated Neumann series
- Three-nested loop involving HVP, i.e., $\nabla^2 f_{i,j}(\mathbf{x}^{(k)}) \mathbf{v}_{i,j}^{(k)}$
- Leverage fast multiplication by Hessian of GLMs

NAME	COMPLEXITY PER ITERATION	REFERENCE
NEWTON-CG	$\mathcal{O}(\text{NNZ}(\mathbf{A})\sqrt{\kappa})$	FOLKLORE
SSN-LS	$\tilde{\mathcal{O}}(\text{NNZ}(\mathbf{A}) \log n + p^2 \kappa^{3/2})$	[XU ET AL., 2016]
SSN-RNS	$\tilde{\mathcal{O}}(\text{NNZ}(\mathbf{A}) + \text{sr}(\mathbf{A})p\kappa^{5/2})$	[XU ET AL., 2016]
SRHT	$\tilde{\mathcal{O}}(np(\log n)^4 + p^2(\log n)^4 \kappa^{3/2})$	[PILANCI ET AL., 2016]
SSN-UNIFORM	$\tilde{\mathcal{O}}(\text{NNZ}(\mathbf{A}) + p\hat{\kappa}\kappa^{3/2})$	[ROOSTA ET AL., 2016]
LISSA	$\tilde{\mathcal{O}}(\text{NNZ}(\mathbf{A}) + p\hat{\kappa}\bar{\kappa}^2)$	[AGARWAL ET AL., 2017]

$$\left. \begin{aligned}
 \kappa &= \max_{\mathbf{x}} \frac{\lambda_{\max} \nabla^2 F(\mathbf{x})}{\lambda_{\min} \nabla^2 F(\mathbf{x})} \\
 \hat{\kappa} &= \max_{\mathbf{x}} \frac{\max_j \lambda_{\max} \nabla^2 f_j(\mathbf{x})}{\lambda_{\min} \nabla^2 F(\mathbf{x})} \\
 \bar{\kappa} &= \max_{\mathbf{x}} \frac{\max_j \lambda_{\max} \nabla^2 f_j(\mathbf{x})}{\min_j \lambda_{\min} \nabla^2 f_j(\mathbf{x})}
 \end{aligned} \right\} \Rightarrow \kappa \leq \hat{\kappa} \leq \bar{\kappa}$$

Sub-sampled Newton's Method

Iterative Scheme

$$\mathbf{p}^{(k)} \approx \operatorname{argmin}_{\mathbf{p} \in \mathcal{R}^d} \left\{ \mathbf{p}^T \mathbf{g}^{(k)} + \frac{1}{2} \mathbf{p}^T \mathbf{H}^{(k)} \mathbf{p} \right\},$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha^{(k)} \mathbf{p}^{(k)}$$

Gradient and Hessian Sub-Sampling

$$\mathbf{g}^{(k)} = \frac{1}{|\mathcal{S}_g|} \sum_{j \in \mathcal{S}_g} \nabla f_j(\mathbf{x}^{(k)})$$

$$\mathbf{H}^{(k)} = \frac{1}{|\mathcal{S}_H|} \sum_{j \in \mathcal{S}_H} \nabla^2 f_j(\mathbf{x}^{(k)})$$

Sub-sampled Newton's Method

Algorithm Newton's Method with **Hessian** and **Gradient** Sub-Sampling

- 1: **Input:** $\mathbf{x}^{(0)}$
 - 2: **for** $k = 0, 1, 2, \dots$ until termination **do**
 - 3: - $\mathcal{S}_g^{(k)} \subseteq \{1, 2, \dots, n\} \implies \mathbf{g}^{(k)} = \frac{1}{|\mathcal{S}_g^{(k)}|} \sum_{j \in \mathcal{S}_g^{(k)}} \nabla f_j(\mathbf{x}^{(k)})$
 - 4: **if** $\|\mathbf{g}^{(k)}\| \leq \sigma \epsilon_g$ **then**
 - 5: - STOP
 - 6: **end if**
 - 7: - $\mathcal{S}_H^{(k)} \subseteq \{1, 2, \dots, n\} \implies \mathbf{H}^{(k)} = \frac{1}{|\mathcal{S}_H^{(k)}|} \sum_{j \in \mathcal{S}_H^{(k)}} \nabla^2 f_j(\mathbf{x}^{(k)})$
 - 8: - $\mathbf{H}^{(k)} \mathbf{p}^{(k)} \approx -\mathbf{g}^{(k)}$
 - 9: - Find $\alpha^{(k)}$ that passes Armijo linesearch
 - 10: - Update $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha^{(k)} \mathbf{p}^{(k)}$
 - 11: **end for**
-

Sub-sampled Newton's Method

Theorem ([Roosta and Mahoney, 2016a])

If

$$\theta \leq \sqrt{\frac{(1 - \epsilon_{\mathbf{H}})}{\kappa}},$$

then, w.h.p,

$$F(\mathbf{x}^{(k+1)}) - F(\mathbf{x}^*) \leq \rho \left(F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*) \right),$$

where $\rho = 1 - (1 - \epsilon_{\mathbf{H}})/\kappa^2$, and upon "STOP", we have $\|\nabla F(\mathbf{x}^{(k)})\| < (1 + \sigma)\epsilon_{\mathbf{g}}$.

Sub-sampled Newton's Method

Theorem ([Roosta and Mahoney, 2016b])

Consider any $0 < \rho_0 + \rho_1 < \rho < 1$. If $\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq c(\rho_0, \rho_1, \rho)$,
 $\epsilon_{\mathbf{g}}^{(k)} = \rho^k \epsilon_{\mathbf{g}}$ and

$$\theta \leq \rho_0 \sqrt{\frac{(1 - \epsilon_{\mathbf{H}})}{\kappa}},$$

we get locally *R-linear* convergence

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq c\rho^k$$

with probability $(1 - \delta)^{2k}$.

- **Problem-independent** local convergence rate

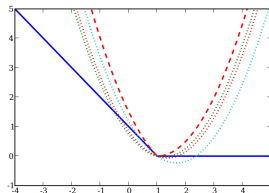
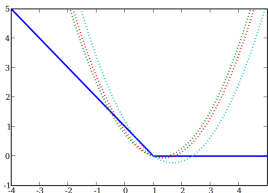
Non-Smooth Newton-type [Yu et al., 2010]

Let $\mathcal{X} = \mathcal{R}^d$, i.e., unconstrained optimization

Iterative Scheme (Non-Quadratic Sub-Problem)

$$\mathbf{p}^{(k)} \approx \operatorname{argmin}_{\mathbf{p} \in \mathcal{R}^d} \left\{ \sup_{\mathbf{g}^{(k)} \in \partial(F+R)(\mathbf{x}^{(k)})} \left\{ \mathbf{p}^T \mathbf{g}^{(k)} \right\} + \frac{1}{2} \mathbf{p}^T \underbrace{\mathbf{H}^{(k)}}_{\substack{\text{e.g.,} \\ \text{QN}}} \mathbf{p} \right\},$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha^{(k)} \mathbf{p}^{(k)}$$



Sub-Sampled Proximal Newton-type [Liu et al., 2017]

FSM/ERM

$$\min_{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d} F(\mathbf{x}) + R(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{a}_i^T \mathbf{x}) + R(\mathbf{x})$$

- f_i : Strongly-Convex, Smooth, and Self-Concordant
- R : Convex and Non-Smooth
- $n \gg 1$ and/or $d \gg 1$
- Dennis-Moré condition:

$$|\mathbf{p}^{(k)T} \left(\mathbf{H}^{(k)} - \nabla^2 F(\mathbf{x}^{(k)}) \right) \mathbf{p}^{(k)}| \leq \eta_k \mathbf{p}^{(k)T} \nabla^2 F(\mathbf{x}^{(k)}) \mathbf{p}^{(k)}$$

- Leverage score sampling to ensure Dennis-Moré
- Inexact sub-problem solver

Finite Sum / Empirical Risk Minimization

Self-Concordant

$$|\mathbf{v}^T (\nabla^3 f(\mathbf{x})[\mathbf{v}]) \mathbf{v}| \leq M (\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v})^{3/2}, \quad \forall \mathbf{x}, \mathbf{v} \in \mathcal{R}^d$$

$M = 2$ is called standard self-concordant.

Theorem ([Zhang and Lin, 2015])

Suppose there exists $\gamma > 0$ and $\eta \in [0, 1)$ such that

$$|f_i'''(t)| \leq \gamma (f_i''(t))^{1-\eta}. \quad \text{Then}$$

$$F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{a}_i^T \mathbf{x}) + \frac{\lambda}{2} \|\mathbf{x}\|^2$$

is self-concordant with

$$M = \frac{\max_i \|\mathbf{a}_i\|^{1+2\eta} \gamma}{\lambda^{\eta+1/2}}.$$

