

# SUB-SAMPLED NEWTON METHODS

F. Roosta-Khorasani and M. W. Mahoney

ICSI and Dept of Statistics, UC Berkeley  
February 2016

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

## PROBLEM STATEMENT

## PROBLEM

$$\min_{\mathbf{x} \in \mathcal{D} \cap \mathcal{X}} F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

- Convex constraint set:  $\mathcal{X} \subseteq \mathbb{R}^p$
- Convex and open domain of  $F$ :  $\mathcal{D} = \bigcap_{i=1}^n \text{dom}(f_i)$
- *Large Scale*:
  - $n \gg 1$
  - $p \gg 1$

# EXAMPLES

- Machine Learning and Data fitting
  - Each  $f_i$  corresponds to an observation (or a measurement) which models the loss (or misfit)
  - Logistic regression
  - SVM
  - Neural Networks
  - Graphical Models
- Nonlinear inverse problems
  - e.g., PDE inverse problems

## ITERATIVE SCHEME

$$x_{k+1} = \arg \min_{\mathbf{x} \in \mathcal{D} \cap \mathcal{X}} \left\{ F(\mathbf{x}^{(k)}) + (\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{g}(\mathbf{x}^{(k)}) + \frac{1}{2\alpha_k} (\mathbf{x} - \mathbf{x}^{(k)})^T H(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) \right\},$$

where

- $\mathbf{g}(\mathbf{x}) \approx \nabla F(\mathbf{x})$
- $H(\mathbf{x}) \approx \nabla^2 F(\mathbf{x})$

- **Newton:**  $\mathbf{g}(\mathbf{x}^{(k)}) = \nabla F(\mathbf{x}^{(k)})$  and  $H(\mathbf{x}^{(k)}) = \nabla^2 F(\mathbf{x}^{(k)})$
- **Projected Gradient Descent:**  $\mathbf{g}(\mathbf{x}^{(k)}) = \nabla F(\mathbf{x}^{(k)})$  and  $H(\mathbf{x}^{(k)}) = \mathbb{I}$
- **Frank-Wolfe:**  $\mathbf{g}(\mathbf{x}^{(k)}) = \nabla F(\mathbf{x}^{(k)})$  and  $H(\mathbf{x}^{(k)}) = 0$
- **(mini-batch) SGD:**  $\mathbf{g}(\mathbf{x}^{(k)}) = 1/|\mathcal{S}_g| \sum_{j \in \mathcal{S}_g} \nabla f_j(\mathbf{x}^{(k)})$  and  $H(\mathbf{x}^{(k)}) = \mathbb{I}$ ,
- **SSN:**
  - **SSN w. Hessian Sub-Sampling:**

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

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

- **SNN w. Gradient and Hessian Sub-Sampling:**

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

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

## MODERN “BIG-DATA”

- $n \gg 1$ 
  - Evaluation  $\nabla F(\mathbf{x})$  and  $\nabla^2 F(\mathbf{x})$  scales “linearly” in  $n$
- $p \gg 1$ 
  - Evaluation of each  $\nabla f_i(\mathbf{x})$  and  $\nabla^2 f_i(\mathbf{x})$  can be expensive
  - The “best” direction of descent  $\Rightarrow$  computationally expensive
- Classical *deterministic* optimization algorithms  $\rightarrow$  **Inefficient**
- Need to design *stochastic* variants
  - Should be **efficient**
  - Should **preserve** as much of original **speed** as possible



- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

## FIRST ORDER METHODS

- Use only gradient information
  - E.g. : Gradient Descent:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla F(\mathbf{x}^{(k)})$$

- Smooth Convex  $F \Rightarrow$  Sublinear,  $\mathcal{O}(1/k)$
- Smooth Strongly Convex  $F \Rightarrow$  Linear,  $\mathcal{O}(\rho^k)$ ,  $\rho < 1$
- However, **iteration cost** scales **linearly** in  $n$

## FIRST ORDER METHODS

- **Stochastic** variants e.g., (mini-batch) SGD
  - $\mathcal{S} \subset \{1, 2, \dots, n\}$  is chosen at random with  $|\mathcal{S}| \ll n$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \sum_{j \in \mathcal{S}} \nabla f_j(\mathbf{x}^{(k)})$$

- **Cheap Per-Iteration costs!**
- However **slower** to converge:
  - Smooth Convex  $F \Rightarrow \mathcal{O}(1/\sqrt{k})$
  - Smooth Strongly Convex  $F \Rightarrow \mathcal{O}(1/k)$
  - Devise **modifications** to
    - **achieve** the convergence **rate** of the **full GD**
    - **preserve** the **per-iteration cost** of **SGD**
    - E.g.: SAG, SDCA, SVRG,...

## SECOND ORDER METHODS

- Use both gradient and Hessian information
  - E.g. : Newton's method:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - [\nabla^2 F(\mathbf{x}^{(k)})]^{-1} \nabla F(\mathbf{x}^{(k)})$$

- Smooth Convex  $F \Rightarrow$  **Locally Superlinear**
- Smooth Strongly Convex  $F \Rightarrow$  **Locally Quadratic**
- However, **iteration cost is high!**

## SECOND ORDER METHODS

- **Approximating** second order information **cheaply**
  - Quasi-Newton, e.g., BFGS and L-BFGS [Nocedal, 1980]
  - Sketching the Hessian [Pilanci et al., 2015]
  - Sub-Sampling the Hessian [Byrd et al., 2011, Erdogdu et al., 2015, Martens, 2010, RM-I & RM-II, 2016]

- $\mathcal{S} \subset \{1, 2, \dots, n\}$  is chosen at random with  $|\mathcal{S}| \ll n$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \left[ \sum_{j \in \mathcal{S}} \nabla^2 f_j(\mathbf{x}^{(k)}) \right]^{-1} \nabla F(\mathbf{x}^{(k)})$$

- Sampling the Hessian and the gradient [RM-I & RM-II, 2016]
  - $\mathcal{S}_H, \mathcal{S}_g \subset \{1, 2, \dots, n\}$  is chosen at random with  $|\mathcal{S}_H|, |\mathcal{S}_g| \ll n$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \left[ \sum_{j \in \mathcal{S}_H} \nabla^2 f_j(\mathbf{x}^{(k)}) \right]^{-1} \sum_{j \in \mathcal{S}_g} \nabla f_j(\mathbf{x}^{(k)})$$

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- **SSN:**
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

## SSN [RM-I &amp; RM-II, 2016]

- Globally Convergent Algorithms [RM-I, 2016]
  - Approach the optimum,  $\mathbf{x}^*$ , from **any**  $\mathbf{x}^{(0)}$
- Local Convergence Rate [RM-II, 2016]
  - Achieve **fast** rate, at least **locally**

Combine  $\Rightarrow$  **Globally convergent** algorithms with **fast local rates!!!**

# ASSUMPTIONS

- **Unconstrained**, i.e.,  $\mathcal{X} = \mathcal{D} = \mathbb{R}^p$
- Each  $f_i$  is **smooth** and **convex**

$$\nabla^2 f_i(\mathbf{x}) \succeq 0, \quad \forall \mathbf{x} \in \mathbb{R}^p, \quad i = 1, 2, \dots, n,$$

$$\nabla^2 f_i(\mathbf{x}) \leq K < \infty, \quad \forall \mathbf{x} \in \mathbb{R}^p, \quad i = 1, 2, \dots, n,$$

$$\|\nabla^2 f_i(\mathbf{x}) - \nabla^2 f_i(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p, \quad i = 1, 2, \dots, n.$$

- $F$  is **strongly convex**

$$\nabla^2 F(\mathbf{x}) \geq \gamma, \quad \forall \mathbf{x} \in \mathbb{R}^p.$$

- condition number:  $\kappa := K/\gamma$

See [RM-II,2016] for **general constraints** and **relaxed assumptions**.



# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

# GLOBALLY CONVERGENT ALGORITHMS

- Requirements:

- (R.1)  $|S|$  must be **independent** of  $n$ , or at least **smaller** than  $n$
- (R.2)  $H(\mathbf{x})$  must be, at least, **invertible**
- (R.3)  $\mathbf{g}(\mathbf{x})$  must be **close** to  $\nabla F(\mathbf{x})$
- (R.4) **Global** convergence guarantee
- (R.5) For  $p \gg 1$ , allow for **inexactness** in solving the linear system

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling  $\Rightarrow$  SSN-H
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

## HESSIAN SUB-SAMPLING

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

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

## SUB-SAMPLING HESSIAN

- Satisfying Requirements (R.1) and (R.2)

## LEMMA (UNIFORM HESSIAN SUB-SAMPLING)

Given any  $0 < \epsilon < 1$ ,  $0 < \delta < 1$ , and  $\mathbf{x} \in \mathbb{R}^p$ , if

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

then

$$\Pr \left( (1 - \epsilon)\gamma \leq \lambda_{\min}(H(\mathbf{x})) \right) \geq 1 - \delta.$$

## SSN-H WITH EXACT UPDATE

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}_k,$$

where

$$\mathbf{p}_k = -[H(\mathbf{x}^{(k)})]^{-1} \nabla F(\mathbf{x}^{(k)}), \quad \text{“exact solve”}$$

$$\alpha_k = \arg \max \quad \alpha$$

$$\text{s.t.} \quad \alpha \leq \hat{\alpha}$$

$$F(\mathbf{x}^{(k)} + \alpha \mathbf{p}_k) \leq F(\mathbf{x}^{(k)}) + \alpha \beta \mathbf{p}_k^T \nabla F(\mathbf{x}^{(k)})$$

$$0 < \beta < 1, \quad \hat{\alpha} \geq 1$$

## SSN-H ALGORITHM: EXACT UPDATE

---

**Algorithm 1** Globally Convergent SSN-H with exact solve

---

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon < 1$ ,  $0 < \beta < 1$ ,  $\hat{\alpha} \geq 1$
  - 2: - Set the sample size,  $|\mathcal{S}|$ , with  $\epsilon$  and  $\delta$
  - 3: **for**  $k = 0, 1, 2, \dots$  until termination **do**
  - 4:   - Select a sample set,  $\mathcal{S}$ , of size  $|\mathcal{S}|$
  - 5:   - Form  $H(\mathbf{x}^{(k)})$
  - 6:   - Update  $\mathbf{x}^{(k+1)}$  with **exact** solve
  - 7: **end for**
-

## GLOBAL CONVERGENCE OF SSN-H: EXACT UPDATE

- Satisfying requirement (R.3)

## THEOREM (GLOBAL CONVERGENCE OF ALGORITHM 1)

Using Algorithm 1 with any  $\mathbf{x}^{(k)} \in \mathbb{R}^p$ , with probability  $1 - \delta$ , we have

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

where  $\rho = 2\alpha_k\beta/\kappa$ . Moreover, the step size is at least

$$\alpha_k \geq 2(1 - \beta)(1 - \epsilon)/\kappa.$$



SSN-H WITH **INEXACT** UPDATE

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}_k,$$

where

$$\|H(\mathbf{x}^{(k)})\mathbf{p}_k + \nabla F(\mathbf{x}^{(k)})\| \leq \theta_1 \|\nabla F(\mathbf{x}^{(k)})\|$$

$$\mathbf{p}_k^T \nabla F(\mathbf{x}^{(k)}) \leq -(1 - \theta_2) \mathbf{p}_k^T H(\mathbf{x}^{(k)}) \mathbf{p}_k$$

$$\alpha_k = \arg \max \quad \alpha$$

$$\text{s.t.} \quad \alpha \leq \hat{\alpha}$$

$$F(\mathbf{x}^{(k)} + \alpha \mathbf{p}_k) \leq F(\mathbf{x}^{(k)}) + \alpha \beta \mathbf{p}_k^T \nabla F(\mathbf{x}^{(k)})$$

$$0 < \beta < 1, \hat{\alpha} \geq 1, 0 < \theta_1, \theta_2 < 1$$

SSN-H ALGORITHM: **INEXACT** UPDATE

---

**Algorithm 2** Globally Convergent SSN-H with inexact solve

---

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon < 1$ ,  $0 < \beta < 1$ ,  $\hat{\alpha} \geq 1$ ,  $0 < \theta_1, \theta_2 < 1$
  - 2: - Set the sample size,  $|\mathcal{S}|$ , with  $\epsilon$  and  $\delta$
  - 3: **for**  $k = 0, 1, 2, \dots$  until termination **do**
  - 4:   - Select a sample set,  $\mathcal{S}$ , of size  $|\mathcal{S}|$
  - 5:   - Form  $H(\mathbf{x}^{(k)})$
  - 6:   - Update  $\mathbf{x}^{(k+1)}$  with **inexact** solve
  - 7: **end for**
-

GLOBAL CONVERGENCE SSN-H: **INEXACT** UPDATE

- Satisfying requirements **(R.3)** & **(R.4)**

## THEOREM (GLOBAL CONVERGENCE OF ALGORITHM 2)

Using Algorithm 2 with any  $\mathbf{x}^{(k)} \in \mathbb{R}^p$ , with probability  $1 - \delta$ , we have

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

where

- (I) if  $\theta_1 \leq \sqrt{\frac{(1-\epsilon)}{4\kappa}}$ , then  $\rho = \alpha_k \beta / \kappa$ ,
- (II) otherwise  $\rho = 2(1 - \theta_2)(1 - \theta_1)^2(1 - \epsilon)\alpha_k \beta / \kappa^2$ .

Moreover, for both cases, the step size is at least  $\alpha_k \geq \frac{2(1-\theta_2)(1-\beta)(1-\epsilon)}{\kappa}$ .

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling  $\Rightarrow$  SSN-GH
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

## GRADIENT &amp; HESSIAN SUB-SAMPLING

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

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

## RANDNLA

$$\nabla F(\mathbf{x}) = \left( \begin{array}{c|c|c|c} \nabla f_1(\mathbf{x}) & \nabla f_2(\mathbf{x}) & \cdots & \nabla f_n(\mathbf{x}) \end{array} \right) \begin{pmatrix} 1/n \\ 1/n \\ \vdots \\ 1/n \end{pmatrix}$$

GRADIENT SUB-SAMPLING: **RANDNLA**

## LEMMA (UNIFORM GRADIENT SUB-SAMPLING)

For a given  $\mathbf{x} \in \mathbb{R}^p$ , let

$$\|\nabla f_i(\mathbf{x})\| \leq G(\mathbf{x}), \quad i = 1, 2, \dots, n.$$

For any  $0 < \epsilon < 1$  and  $0 < \delta < 1$ , if

$$|\mathcal{S}| \geq \frac{G(\mathbf{x})^2}{\epsilon^2} \left(1 + \sqrt{8 \ln \frac{1}{\delta}}\right)^2,$$

then

$$\Pr \left( \|\nabla F(\mathbf{x}) - \mathbf{g}(\mathbf{x})\| \leq \epsilon \right) \geq 1 - \delta.$$

- Need to efficiently estimate  $G(\mathbf{x})$  at every iteration...see examples in [\[RM-I,2016\]](#)

## SSN-GH WITH EXACT UPDATE

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}_k,$$

where

$$\mathbf{p}_k = -[H(\mathbf{x}^{(k)})]^{-1} \mathbf{g}(\mathbf{x}^{(k)}), \quad \text{“exact solve”}$$

$$\alpha_k = \arg \max \quad \alpha$$

$$\text{s.t.} \quad \alpha \leq \hat{\alpha}$$

$$F(\mathbf{x}^{(k)} + \alpha \mathbf{p}_k) \leq F(\mathbf{x}^{(k)}) + \alpha \beta \mathbf{p}_k^T \mathbf{g}(\mathbf{x}^{(k)})$$

$$0 < \beta < 1, \hat{\alpha} \geq 1$$



## SSN-GH WITH EXACT UPDATE

**Algorithm 3** Globally Convergent SSN-GH with exact solve

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon_1 < 1$ ,  $0 < \epsilon_2 < 1$ ,  $0 < \beta < 1$ ,  $\hat{\alpha} \geq 1$  and  $\sigma \geq 0$
- 2: - Set the sample size,  $|\mathcal{S}_H|$ , with  $\epsilon_1$  and  $\delta$
- 3: **for**  $k = 0, 1, 2, \dots$  until termination **do**
- 4:   - Select a sample set,  $\mathcal{S}_H$ , of size  $|\mathcal{S}_H|$  and form  $H(\mathbf{x}^{(k)})$
- 5:   - Set the sample size,  $|\mathcal{S}_g|$ , with  $\epsilon_2$ ,  $\delta$  and  $\mathbf{x}^{(k)}$
- 6:   - Select a sample set,  $\mathcal{S}_g$  of size  $|\mathcal{S}_g|$  and form  $\mathbf{g}(\mathbf{x}^{(k)})$
- 7:   **if**  $\|\mathbf{g}(\mathbf{x}^{(k)})\| < \sigma\epsilon_2$  **then**
- 8:     - STOP
- 9:   **end if**
- 10:   - Update  $\mathbf{x}^{(k+1)}$  with **exact** solve
- 11: **end for**

## GLOBAL CONVERGENCE SSN-GH: EXACT UPDATE

## THEOREM (GLOBAL CONVERGENCE OF ALGORITHM 3)

Using Algorithm 3 with any  $\mathbf{x}^{(k)} \in \mathbb{R}^p$ ,  $\epsilon_1 \leq 1/2$  and  $\sigma \geq 4\kappa/(1 - \beta)$ , we have the following with probability  $(1 - \delta)^2$ :

(I) if “STOP”, then

$$\|\nabla F(\mathbf{x}^{(k)})\| < (1 + \sigma)\epsilon_2,$$

(II) otherwise, we have

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

with  $\rho = 8\alpha_k\beta/(9\kappa)$  with the step size of at least

$$\alpha_k \geq (1 - \beta)(1 - \epsilon_1)/\kappa.$$

## SSN-GH WITH INEXACT UPDATE

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}_k,$$

where

$$\|H(\mathbf{x}^{(k)})\mathbf{p}_k + \mathbf{g}(\mathbf{x}^{(k)})\| \leq \theta_1 \|\mathbf{g}(\mathbf{x}^{(k)})\|$$

$$\mathbf{p}_k^T \mathbf{g}(\mathbf{x}^{(k)}) \leq -(1 - \theta_2) \mathbf{p}_k^T H(\mathbf{x}^{(k)}) \mathbf{p}_k$$

$$\alpha_k = \arg \max \alpha$$

$$\text{s.t. } \alpha \leq \hat{\alpha}$$

$$F(\mathbf{x}^{(k)} + \alpha \mathbf{p}_k) \leq F(\mathbf{x}^{(k)}) + \alpha \beta \mathbf{p}_k^T \mathbf{g}(\mathbf{x}^{(k)})$$

$$0 < \beta < 1, \hat{\alpha} \geq 1, 0 < \theta_1, \theta_2 < 1$$

SSN-GH WITH **INEXACT** UPDATE

---

**Algorithm 4** Globally Convergent SSN-GH with exact solve

---

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon_1 < 1$ ,  $0 < \epsilon_2 < 1$ ,  $0 < \beta < 1$ ,  $\hat{\alpha} \geq 1$ ,  
 $\sigma \geq 0$ ,  $0 < \theta_1, \theta_2 < 1$
  - 2: - Set the sample size,  $|\mathcal{S}_H|$ , with  $\epsilon_1$  and  $\delta$
  - 3: **for**  $k = 0, 1, 2, \dots$  until termination **do**
  - 4:   - Select a sample set,  $\mathcal{S}_H$ , of size  $|\mathcal{S}_H|$  and form  $H(\mathbf{x}^{(k)})$
  - 5:   - Set the sample size,  $|\mathcal{S}_g|$ , with  $\epsilon_2$ ,  $\delta$  and  $\mathbf{x}^{(k)}$
  - 6:   - Select a sample set,  $\mathcal{S}_g$  of size  $|\mathcal{S}_g|$  and form  $\mathbf{g}(\mathbf{x}^{(k)})$
  - 7:   **if**  $\|\mathbf{g}(\mathbf{x}^{(k)})\| < \sigma\epsilon_2$  **then**
  - 8:     - STOP
  - 9:   **end if**
  - 10:   - Update  $\mathbf{x}^{(k+1)}$  with **inexact** solve
  - 11: **end for**
-

# GLOABL CONVERGENCE OF SSN-GH: **INEXACT** UPDATE

## THEOREM (GLOBAL CONVERGENCE OF ALGORITHM 4)

Using Algorithm 4 with  $\epsilon_1 \leq 1/2$ , any  $\mathbf{x}^{(k)} \in \mathbb{R}^p$ , and  $\sigma \geq \frac{4\kappa}{(1-\theta_1)(1-\theta_2)(1-\beta)}$ , we have the following with probability  $1 - \delta$ :

- (I) if “STOP”, then  $\|\nabla F(\mathbf{x}^{(k)})\| < (1 + \sigma)\epsilon_2$ ,
- (II) otherwise  $F(\mathbf{x}^{(k+1)}) - F(\mathbf{x}^*) \leq (1 - \rho)(F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*))$ , where
  - (1) if  $\theta_1 \leq \sqrt{(1 - \epsilon_1)/(4\kappa)}$ , then  $\rho = 4\alpha_k\beta/(9\kappa)$ ,
  - (2) otherwise  $\rho = 8\alpha_k\beta(1 - \theta_2)(1 - \theta_1)^2(1 - \epsilon_1)/(9\kappa^2)$ .

Moreover, for both cases, the step size is at least

$$\alpha_k \geq (1 - \theta_2)(1 - \beta)(1 - \epsilon_1)/\kappa.$$

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

# LOCAL CONVERGENCE RATES

- Requirements:

(R.1)  $|\mathcal{S}|$  must be **independent** of  $n$ , or at least **smaller** than  $n$

(R.2)  $H(\mathbf{x})$  must **preserve** the spectrum of  $\nabla^2 F(\mathbf{x})$  as much as possible

(R.3)  $\mathbf{g}(\mathbf{x})$  must be **close** to  $\nabla F(\mathbf{x})$

(R.4) **Fast** local convergence rate, close to that of Newton!

(R.5) For  $p \gg 1$ , allow for inexactness  $\Rightarrow$  future work

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples



## HESSIAN SUB-SAMPLING

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

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

## SUB-SAMPLING HESSIAN

- Satisfying Requirements (R.1) and (R.2)

## LEMMA (UNIFORM HESSIAN SUB-SAMPLING)

Given any  $0 < \epsilon < 1$ ,  $0 < \delta < 1$  and  $\mathbf{x} \in \mathbb{R}^p$ , if

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

then

$$\Pr \left( |\lambda_i(\nabla^2 F(\mathbf{x})) - \lambda_i(H(\mathbf{x}))| \leq \epsilon \lambda_i(\nabla^2 F(\mathbf{x})); i = 1, 2, \dots, p \right) \geq 1 - \delta.$$

- Difference between (R.2) here and (R.2) before  $\Rightarrow \kappa^2$  here vs.  $\kappa$  before!

## ERROR RECURSION: HESSIAN SUB-SAMPLING

## THEOREM (ERROR RECURSION)

Let  $0 < \delta < 1$  and  $0 < \epsilon < 1$  be given. Using  $\alpha_k = 1$ , with probability  $1 - \delta$ , we have

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \rho_0 \|\mathbf{x}^{(k)} - \mathbf{x}^*\| + \xi \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2,$$

where

$$\rho_0 = \frac{\epsilon}{(1 - \epsilon)}, \quad \text{and} \quad \xi = \frac{L}{2(1 - \epsilon)\gamma}.$$

- $\rho_0$  is **problem-independent!**  $\Rightarrow$  Can be made **arbitrarily small!**

## SSN-H ALGORITHM

---

**Algorithm 5** Locally Convergent SSN-H with exact solve

---

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon < 1$
  - 2: - Set the sample size,  $|\mathcal{S}|$ , with  $\epsilon$  and  $\delta$
  - 3: **for**  $k = 0, 1, 2, \dots$  until termination **do**
  - 4:   - Select a sample set,  $\mathcal{S}$ , of size  $|\mathcal{S}|$  and  $H(\mathbf{x}^{(k)})$
  - 5:   - Update  $\mathbf{x}^{(k+1)}$  with  $H(\mathbf{x}^{(k)})$  and  $\alpha_k = 1$
  - 6: **end for**
-

## SSN-H: Q-LINEAR CONVERGENCE

## THEOREM (Q-LINEAR CONVERGENCE)

Consider any  $0 < \rho_0 < \rho < 1$  and  $\epsilon \leq \rho_0/(1 + \rho_0)$ . If

$$\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq \frac{\rho - \rho_0}{\xi},$$

using Algorithm 5, 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  $(1 - \delta)^{k_0}$ .

## SSN-H: Q-SUPERLINEAR CONVERGENCE

## THEOREM (Q-SUPERLINEAR CONVERGENCE: GEOMETRIC GROWTH)

Using Algorithm 5, with

$$\epsilon^{(k)} = \rho^k \epsilon, \quad k = 0, 1, \dots, k_0,$$

if  $\mathbf{x}^{(0)}$  is close-enough, we get locally *Q-superlinear* convergence

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

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

# HESSIAN SUB-SAMPLING: Q-SUPERLINEAR CONVERGENCE

## THEOREM (Q-SUPERLINEAR CONVERGENCE: SLOW GROWTH)

Using Algorithm 5 with

$$\epsilon^{(k)} = \frac{1}{1 + 2 \ln(4 + k)}, \quad k = 0, 1, \dots, k_0,$$

if  $\mathbf{x}^{(0)}$  is close-enough, we get locally *Q-superlinear* convergence

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \frac{1}{\ln(3 + k)} \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|, \quad k = 1, \dots, k_0,$$

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

# LOCAL VS. GLOBAL

How do we **connect** the **global** and **local** results?



# LOCAL + GLOBAL

## THEOREM (GLOBAL CONV. OF ALG. 1 WITH PROBLEM-INDEPENDENT LOCAL RATE)

Consider any  $0 < \rho_0 < \rho_1 < 1/\sqrt{\kappa}$ . Using Algorithm 1 with any  $\mathbf{x}^{(0)} \in \mathbb{R}^P$ ,  $\hat{\alpha} = 1$ ,  $\epsilon \leq \rho_0 / ((1 + \rho_0)\sqrt{2\kappa})$  and  $\beta \leq (1 - \epsilon)(1 - \kappa\rho_1^2)/(2\kappa)$ , after

$$k \geq 2 \ln \left( \frac{2(\rho_1 - \rho_0)(1 - \epsilon)\gamma}{\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \sqrt{\kappa} L} \right) / \ln(1 - 2\beta/\kappa)$$

iterations, with probability  $(1 - \delta)^k$  we get “*problem-independent*” linear convergence, i.e.,

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \rho_1 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|.$$

Moreover, the step size of  $\alpha_k = 1$  passes Armijo rule for *all* subsequent iterations.

MODIFYING  $H(\mathbf{x})$ 

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \rho_0 \|\mathbf{x}^{(k)} - \mathbf{x}^*\| + \xi \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2$$

$$\rho_0 = \frac{\epsilon}{(1 - \epsilon)} \quad \xi = \frac{L}{2(1 - \epsilon)\gamma}$$

$$\gamma \downarrow \Rightarrow \xi \uparrow$$

## SPECTRAL REGULARIZATION

$$\hat{H} := \mathcal{P}(\lambda; H) = \sum_{i=1}^p \max\{\lambda_i, \lambda\} \mathbf{v}_i \mathbf{v}_i^T$$

- $(\lambda_i, \mathbf{v}_i)$ : eigen-pair of  $H(\mathbf{x})$

## SPECTRAL REGULARIZATION

## THEOREM (ERROR RECURSION)

Let  $0 < \delta < 1$  and  $0 < \epsilon < 1$  be given. Set  $|\mathcal{S}|$  and form  $H(\mathbf{x}^{(k)})$ . For some  $\lambda > 0$ , let

$$\hat{H}(\mathbf{x}^{(k)}) = \mathcal{P}(\lambda; H(\mathbf{x}^{(k)})).$$

Then, with  $\alpha_k = 1$ , if  $\lambda \geq (1 - \epsilon)\gamma$ , it follows that, with probability  $1 - \delta$

$$\rho_0 = \frac{\lambda - (1 - \epsilon)\gamma + \gamma\epsilon}{\lambda}, \quad \text{and} \quad \xi = \frac{L}{2\lambda}.$$

## SPECTRAL REGULARIZATION

**Algorithm 6** SSN-H and Spectral Regularization

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon < 1$ , and  $0 < \epsilon_0 < 1$
- 2: - Set the sample size,  $|\mathcal{S}_0|$ , with  $\epsilon_0$  and  $\delta$
- 3: - Set the sample size,  $|\mathcal{S}|$ , with  $\epsilon$  and  $\delta$
- 4: **for**  $k = 0, 1, 2, \dots$  until termination **do**
- 5:   - Select a sample set,  $\mathcal{S}_0$ , of size  $|\mathcal{S}_0|$  and form  $H_0(\mathbf{x}^{(k)})$
- 6:   - Compute  $\lambda_{\min}(H_0(\mathbf{x}^{(k)}))$
- 7:   - Set the threshold,  $\lambda^{(k)}$
- 8:   - Select a sample set,  $\mathcal{S}$ , of size  $|\mathcal{S}|$  and form  $H(\mathbf{x}^{(k)})$
- 9:   - Form  $\hat{H}(\mathbf{x}^{(k)})$  with  $H(\mathbf{x}^{(k)})$  and  $\lambda^{(k)}$
- 10:   - Update  $\mathbf{x}^{(k+1)}$  with  $\hat{H}(\mathbf{x}^{(k)})$  and  $\alpha_k = 1$
- 11: **end for**

## SPECTRAL REGULARIZATION

## THEOREM (Q-LINEAR CONVERGENCE OF ALGORITHM 6)

Using Algorithm 6, if  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\| < \gamma/(3L)$ ,  $\epsilon \leq 1/6$  and at every iteration the threshold is chosen as

$$\lambda^{(k)} \geq \left( \frac{1 - \epsilon}{1 - \epsilon_0} \right) \lambda_{\min} \left( H_0(\mathbf{x}^{(k)}) \right),$$

we get locally Q-linear convergence with variable rates

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

with probability  $(1 - \delta)^{2k_0}$ , where

$$\rho^{(k)} = 1 - \frac{\gamma}{2\lambda^{(k)}}.$$

# RIDGE REGULARIZATION

$$\hat{H}(\mathbf{x}) := H(\mathbf{x}) + \lambda \mathbb{I}$$

## RIDGE REGULARIZATION

## THEOREM (ERROR RECURSION)

Let  $0 < \delta < 1$  and  $0 < \epsilon < 1$  be given. Set  $|S|$  and form  $H(\mathbf{x}^{(k)})$ . For some  $\lambda \geq 0$ , form  $\hat{H}(\mathbf{x}^{(k)})$ . Then, with  $\alpha_k = 1$ , it follows that, with probability  $1 - \delta$ ,

$$\rho_0 = \frac{\lambda + \gamma\epsilon}{(1 - \epsilon)\gamma + \lambda}, \quad \xi = \frac{L}{2(1 - \epsilon)\gamma + 2\lambda}.$$



# RIDGE REGULARIZATION

---

**Algorithm 7** SSN-H and Ridge Regularization

---

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon < 1$ ,  $\lambda \geq 0$
  - 2: - Set the sample size,  $|\mathcal{S}|$ , with  $\epsilon$  and  $\delta$
  - 3: **for**  $k = 0, 1, 2, \dots$  until termination **do**
  - 4:   - Select a sample set,  $\mathcal{S}$ , of size  $|\mathcal{S}|$  and form  $H(\mathbf{x}^{(k)})$
  - 5:   - Form  $\hat{H}(\mathbf{x}^{(k)})$  with  $H(\mathbf{x}^{(k)})$  and  $\lambda$
  - 6:   - Update  $\mathbf{x}^{(k+1)}$  with  $\hat{H}(\mathbf{x}^{(k)})$  and  $\alpha_k = 1$
  - 7: **end for**
-

# RIDGE REGULARIZATION

## THEOREM (Q-LINEAR CONVERGENCE OF ALGORITHM 7)

For any  $\lambda \geq 0$ , consider  $\rho_0$  and  $\rho$  such that

$$1 - \frac{\gamma}{\gamma + \lambda} < \rho_0 < \rho < 1.$$

Using Algorithm 7 with

$$\epsilon \leq \frac{\rho_0 \gamma + (\rho_0 - 1) \lambda}{(1 + \rho_0) \gamma},$$

if  $\mathbf{x}^{(0)}$  is close enough, then with probability  $(1 - \delta)^{k_0}$ , we get locally Q-linear convergence with the rate  $\rho$ .

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

## GRADIENT &amp; HESSIAN SUB-SAMPLING

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

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

### THEOREM (ERROR RECURSION: INDEPENDENT SUB-SAMPLING)

Let  $0 < \delta < 1$ ,  $0 < \epsilon_1 < 1$ , and  $0 < \epsilon_2 < 1$  be given. Set  $|\mathcal{S}_H|$  with  $(\epsilon_1, \delta)$  and  $|\mathcal{S}_g|$  with  $(\epsilon_2, \delta)$ . Independently, choose  $\mathcal{S}_H$  and  $\mathcal{S}_g$ , and form  $H(\mathbf{x}^{(k)})$  and  $\mathbf{g}(\mathbf{x}^{(k)})$ . Then using  $\alpha_k = 1$ , with probability  $(1 - \delta)^2$ , we have

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \eta + \rho_0 \|\mathbf{x}^{(k)} - \mathbf{x}^*\| + \xi \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2,$$

where

$$\eta = \frac{\epsilon_2}{(1 - \epsilon_1)\gamma}, \quad \rho_0 = \frac{\epsilon_1}{(1 - \epsilon_1)}, \quad \text{and} \quad \xi = \frac{L}{2(1 - \epsilon_1)\gamma}.$$

See [\[RM-II,2016\]](#) for **simultaneous sampling**.

## SSN-GH ALGORITHM

---

**Algorithm 8** Locally Convergent SSN-GH

---

- 1: **Input:**  $\mathbf{x}^{(0)}$ ,  $0 < \delta < 1$ ,  $0 < \epsilon_1 < 1$ ,  $0 < \epsilon_2 < 1$  and  $0 < \rho < 1$
  - 2: - Set the sample size,  $|\mathcal{S}_H|$ , with  $\epsilon_1$  and  $\delta$
  - 3: **for**  $k = 0, 1, 2, \dots$  until termination **do**
  - 4:   - Select a sample set,  $\mathcal{S}_H$ , of size  $|\mathcal{S}_H|$  and form  $H(\mathbf{x}^{(k)})$
  - 5:   - Set  $\epsilon_2^{(k)} = \rho^k \epsilon_2$
  - 6:   - Set the sample size,  $|\mathcal{S}_g^{(k)}|$ , with  $\epsilon_2^{(k)}$ ,  $\delta$  and  $\mathbf{x}^{(k)}$
  - 7:   - Select a sample set,  $\mathcal{S}_g^{(k)}$  of size  $|\mathcal{S}_g^{(k)}|$  and form  $\mathbf{g}(\mathbf{x}^{(k)})$
  - 8:   - Update  $\mathbf{x}^{(k+1)}$  with  $H(\mathbf{x}^{(k)})$ ,  $\mathbf{g}(\mathbf{x}^{(k)})$  and  $\alpha_k = 1$
  - 9: **end for**
-

# R-LINEAR CONVERGENCE OF SSN-GH

## THEOREM (R-LINEAR CONVERGENCE)

Consider any  $0 < \rho < 1$ ,  $0 < \rho_0 < 1$ , and  $0 < \rho_1 < 1$  such that  $\rho_0 + \rho_1 < \rho$ . Let  $\epsilon_1 \leq \rho_0 / (1 + \rho_0)$ , and define  $\sigma := 2(\rho - (\rho_0 + \rho_1))(1 - \epsilon_1)\gamma / L$ . Using Algorithm 8 with  $\epsilon_2 \leq (1 - \epsilon_1)\gamma\rho_1\sigma$ , if the initial iterate,  $\mathbf{x}^{(0)}$ , satisfies  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq \sigma$ , we get locally **R-linear** convergence

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

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

# LOCAL VS. GLOBAL

How do we **connect** the **global** and **local** results?



## THEOREM (GLOBAL CONV. OF ALG. 3 WITH PROBLEM-INDEPENDENT LOCAL RATE)

Consider any  $0 < \rho_0, \rho_1, \rho_2 < 1/\sqrt{\kappa}$  such that  $\rho_0 + \rho_1 < \rho_2$ , set  $\epsilon_1 \leq \rho_0/((1 + \rho_0)\sqrt{2\kappa})$ . Using Algorithm 3 with any  $\mathbf{x}^{(0)} \in \mathbb{R}^p$  and

$$\hat{\alpha} = 1, \quad \beta \leq \frac{(1 - \epsilon_1)(1 - \kappa\rho_2^2)}{8\kappa}$$

$$\epsilon_2^{(0)} \leq (1 - \epsilon_1)\gamma\rho_1\sigma, \quad \epsilon_2^{(k)} = \rho_2\epsilon_2^{(k-1)}, \quad k = 1, 2, \dots,$$

after

$$k \geq 2 \ln \left( \frac{\sigma}{\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \sqrt{\kappa}} \right) / \ln(1 - 8\beta/(9\kappa))$$

iterations, we have the following with probability  $(1 - \delta)^{2k}$ :

- 1 if “STOP”, then  $\|\nabla F(\mathbf{x}^{(k)})\| < (1 + \sigma)\rho_2^k\epsilon_2^{(0)}$
- 2 otherwise, we get “**problem-independent**” linear convergence, i.e.,  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \rho_2\|\mathbf{x}^{(k)} - \mathbf{x}^*\|$ . Moreover, the step size of  $\alpha_k = 1$  passes Armiju rule for **all** subsequent iterations.

# OUTLINE

- Problem Statement
- Rough Overview of Existing Methods
  - First Order methods
  - Second order methods
- SSN:
  - Globally convergent algorithms
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
  - Local convergent rates
    - Hessian Sub-Sampling
    - Gradient & Hessian Sub-Sampling
- Examples

## GLM WITH GAUSSIAN PRIOR

$$F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \left( \Phi(\mathbf{a}_i^T \mathbf{x}) - b_i \mathbf{a}_i^T \mathbf{x} \right) + \frac{\lambda}{2} \|\mathbf{x}\|^2.$$

- $\Phi$ : **cumulant generating function**
  - $\Phi(t) = 0.5t^2 \Rightarrow$  Ridge Regression (**RR**)
  - $\Phi(t) = \ln(1 + \exp(t)) \Rightarrow \ell_2$  regularized Logistic Regression (**LR**)
  - $\Phi(t) = \exp(t) \Rightarrow \ell_2$  regularized Poisson Regression (**PR**)

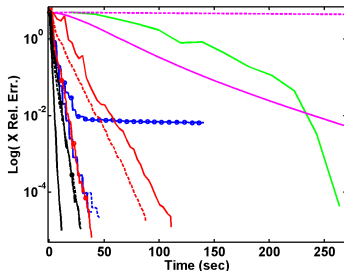
$\ell_2$  REGULARIZED LOGISTIC REGRESSION

	LR
$\nabla f_i(\mathbf{x})$	$\left( \frac{1}{1+e^{-\mathbf{a}_i^T \mathbf{x}}} - b_i \right) \mathbf{a}_i + \lambda \mathbf{x}$
$\nabla^2 f_i(\mathbf{x})$	$\frac{e^{\mathbf{a}_i^T \mathbf{x}}}{(e^{\mathbf{a}_i^T \mathbf{x}} + 1)^2} \mathbf{a}_i \mathbf{a}_i^T + \lambda \mathbb{I}$
$K$	$0.25 \max_i \ \mathbf{a}_i\ ^2 + \lambda$
$\gamma$	$\lambda$
$G(\mathbf{x})$	$\lambda \ \mathbf{x}\  + \max_i (1 +  b_i ) \ \mathbf{a}_i\ $

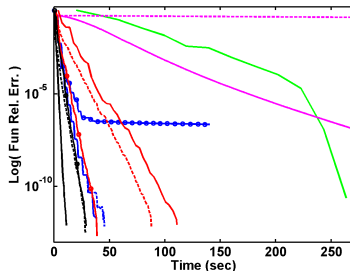
# NUMERICAL EXAMPLES

- $\ell_2$  regularized Logistic Regression

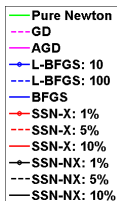
DATA	$n$	$p$	NNZ	$\kappa$
$D_1$	$10^6$	$10^4$	0.02%	$\approx 10^4$
$D_2$	$5 \times 10^4$	$5 \times 10^3$	DENSE	$\approx 10^6$
$D_3$	$10^7$	$2 \times 10^4$	0.006%	$\approx 10^{10}$

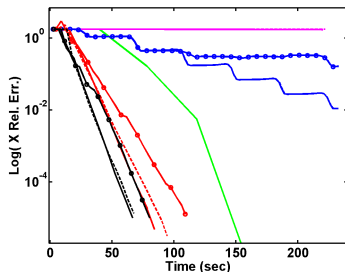
$D_1, n = 10^6, \rho = 10^4, \text{SPARSITY} : 0.02\%, \kappa \approx 10^4$ 


(a) Iterate Relative Error

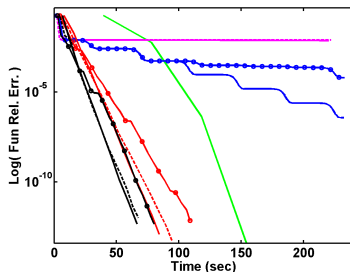


(b) Function Relative Error

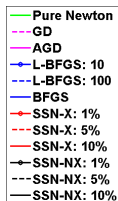


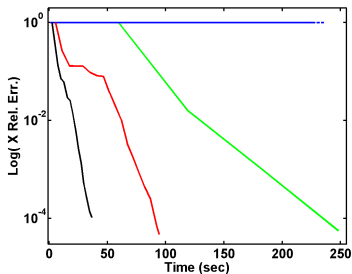
$D_2, n = 5 \times 10^4, p = 5 \times 10^3, \text{SPARSITY : DENSE}, \kappa \approx 10^6$ 


(d) Iterate Relative Error

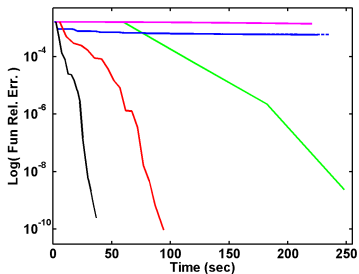


(e) Function Relative Error

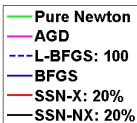


$D_3, n = 10^7, p = 2 \times 10^4, \text{SPARSITY} : 0.006\%, \kappa \approx 10^{10}$ 


(g) Iterate Relative Error



(h) Function Relative Error





THANK YOU!