

Background

Consider the regularized optimization problem,

$$\overline{\theta} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta) + \mathcal{R}_{\lambda_{\operatorname{tgt}}}(\theta), \tag{1}$$

• $\mathcal{L}: \mathbb{R}^d \to \mathbb{R}$ is a twice differentiable convex loss function, e.g., negative log-likelihood for GLM,

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell_i(\theta), \ \ell_i(\theta) = \psi(x_i^{\top}\theta) - y_i x_i^{\top}\theta,$$

• $\mathcal{R}_{\lambda_{tgt}} : \mathbb{R}^d \to \mathbb{R}$ is a sparsity-inducing decomposable regularizer, i.e., $\mathcal{R}_{\lambda_{\text{tgt}}}(\theta) = \sum_{j=1}^{d} r_{\lambda_{\text{tgt}}}(\theta_j)$ with $r_{\lambda_{\text{tgt}}} : \mathbb{R} \to \mathbb{R}$, and $\lambda_{tgt} > 0$ is the regularization parameter.

Convex vs. Nonconvex $\mathcal{R}_{\lambda_{tgt}}$:

- Convex: easy optimization, large estimation bias (e.g., ℓ_1)
- Nonconvex: small estimation bias, hard optimization (e.g., SCAD, MCP, and Capped ℓ_1)



Consider the Capped ℓ_1 regularizer [2] defined as

$$\mathcal{R}_{\lambda_{\text{tgt}}}(\theta) = \sum_{j=1}^{d} r_{\text{tgt}}(\theta_j) = \lambda_{\text{tgt}} \sum_{j=1}^{d} \min\{|\theta_j|, \beta \lambda_{\text{tgt}}\}, \quad (2)$$

where $\beta > 0$ is a tuning parameter. $r_{\lambda_{tgt}}(\theta_j)$ can be decomposed as the **difference of convex functions**:



Our Approach

Algorithm 1 DC Proximal Newton Algorithm

Input: $\widehat{ heta}^{\{0\}}$, $\lambda_{ ext{tgt}}$, eta, arepsilonWarm Initialization: $\hat{\theta}^{\{1\}} \leftarrow \text{ProxNewton}(\hat{\theta}^{\{0\}}, \lambda_{\text{tgt}}, \varepsilon), K \leftarrow 1$ **Repeat:**
$$\begin{split} \lambda_{j}^{\{K+1\}} \leftarrow \begin{cases} 0, & \text{if } |\widehat{\theta}_{j}^{\{K\}}| > \beta \lambda_{\text{tgt}} \\ \lambda_{\text{tgt}}, & \text{if } |\widehat{\theta}_{j}^{\{K\}}| \le \beta \lambda_{\text{tgt}} \\ t \leftarrow 0, \ \theta^{(0)} = \widehat{\theta}^{\{K\}} \end{split}$$
Repeat: $\theta^{(t+1)} \leftarrow \operatorname{argmin}_{\theta} \mathcal{Q}(\theta; \theta^{(t)}, \lambda^{\{K+1\}})$ $t \leftarrow t + 1$ **Until** $\omega_{\lambda\{K+1\}}(\theta^{(t)}) \leq \varepsilon$ $\widehat{\theta}^{\{K+1\}} \leftarrow \theta^{(t)}$ $K \leftarrow K + 1$ **Until** Convergence **Return:** $\hat{\theta}^{\{K\}}$

On Quadratic Convergence of DC Proximal Newton Algorithm in Nonconvex Sparse Learning

Xingguo Li^{*}, Lin Yang^{\$}, Jian Ge^{\$}, Jarvis Haupt^{*}, Tong Zhang^{*}, Tuo Zhao[†] *University of Minnesota *Princeton University *Tencent Lab [†]Georgia Tech

DC Proximal Newton Algorithm

Our DC proximal Newton algorithm contains three components: (I) The multistage convex relaxation. At the $\{K+1\}$ -th stage, we solve a convex relaxation of (1) at $\theta = \hat{\theta}^{\{K\}}$ as

$$\bar{\theta}^{\{K+1\}} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta) + ||\lambda^{\{K+1\}} \odot \theta||_1, \quad (3)$$

where
$$\lambda^{\{K+1\}} = \left(\lambda_1^{\{K+1\}}, ..., \lambda_d^{\{K+1\}}\right)^{\top}, \lambda_j^{\{K+1\}} = \lambda_{\text{tgt}} \cdot \left(|\widehat{\theta}_i^{\{K\}}| \le \beta \lambda_{\text{tgt}}\right)$$
 for all $j = 1, ..., d$.

(II) The warm initialization. In the first stage of DC programming, we solve the ℓ_1 -regularized counterpart of (1) as

$$\overline{\theta}^{\{1\}} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta) + \lambda_{\operatorname{tgt}} ||\theta||_1.$$
(4)

(III) The proximal Newton algorithm. At t-th iteration of the proximal Newton algorithm, we consider a quadratic approximation of (3) at $heta^{(t)}$ as

$$\mathcal{Q}(\theta; \theta^{(t)}, \lambda^{\{K\}}) = \mathcal{L}(\theta^{(t)}) + (\theta - \theta^{(t)})^{\top} \nabla \mathcal{L}(\theta^{(t)}) + \frac{1}{2} ||\theta - \theta^{(t)}||_{\nabla^2 \mathcal{L}(\theta^{(t)})}^2 + ||\lambda^{\{K\}} \odot \theta||_1, \quad (5)$$

where
$$\|\theta - \theta^{(t)}\|^2_{\nabla^2 \mathcal{L}(\theta^{(t)})} = (\theta - \theta^{(t)})^\top \nabla^2 \mathcal{L}(\theta^{(t)})(\theta - \theta^{(t)})$$
. Then
 $\theta^{(t+1)} = \operatorname{argmin}_{\theta} \mathcal{Q}(\theta; \theta^{(t)}, \lambda^{\{K\}}).$

Backtracking line search is used in the warm initialization. We terminate the iterations when the *approximate KKT condition* holds:

$$\omega_{\lambda^{\{1\}}}\left(\theta^{(t)}\right) = \min_{\xi \in \partial ||\theta^{(t)}||_{1}} ||\nabla \mathcal{L}(\theta^{(t)}) + \lambda^{\{1\}} \odot \xi||_{\infty} \le \varepsilon,$$

Assumptions

Assumption 1 (Local Restricted Sparse Eigenvalue (SE)). Given $\theta \in \mathcal{B}(\theta^*, R) = \{\phi \in \mathbb{R}^d \mid ||\phi - \theta^*||_2 \leq R\}$ for a generic constant R, \exists a constant C_0 such that $\nabla^2 \mathcal{L}(\theta)$ satisfies the SE properties with parameters $\rho_{s^*+2\widetilde{s}}^-$ and $\rho_{s^*+2\widetilde{s}}^+$ satisfying

$$0 < \rho_{s^*+2\widetilde{s}}^- = \inf_v \frac{v^\top \nabla^2 \mathcal{L}(\theta) v}{v^\top v} < \rho_{s^*+2\widetilde{s}}^+ = \sup_v \frac{v^\top \nabla^2 \mathcal{L}(\theta) v}{v^\top v} < +\infty$$

where $\|v\|_0 \leq s^* + 2\widetilde{s}$, $\widetilde{s} \geq$ $C_0 \kappa_{s^*+2\widetilde{s}}^2 s^*$, and $\kappa_{s^*+2\widetilde{s}} =$ $\rho_{s^*+2\widetilde{s}}^{-}/\rho_{s^*+2\widetilde{s}}^{-}.$



Assumption 2 (Local Restricted Hessian Smoothness). \exists generic constants $L_{s^*+2\widetilde{s}}$ and R such that for any $\theta, \theta' \in \mathcal{B}(\theta^*, R)$ with $||\theta_{\mathcal{S}_{\perp}}||_0 \leq \widetilde{s} \text{ and } ||\theta'_{\mathcal{S}_{\perp}}||_0 \leq \widetilde{s}, \text{ we have }$

 $||\nabla^2 \mathcal{L}(\theta) - \nabla^2 \mathcal{L}(\theta')||_2 \le L_{s^* + 2\widetilde{s}} ||\theta - \theta'||_2.$

Assumption 3. Given the true parameter θ^* , \exists generic constant C_1 such that $\lambda_{ ext{tgt}} = C_1 \sqrt{rac{\log d}{n}} \geq 4 ||
abla \mathcal{L}(heta^*) ||_\infty$. Moreover, for large enough n, we have $\sqrt{s^*}\lambda_{tgt} \leq C_2 R \rho_{s^*+2\widetilde{s}}^-$ with $R = \frac{\rho_{s^*+2\widetilde{s}}}{2L_{s^*+2\widetilde{s}}}$.

Assumption 4. For each stage of solving the convex relaxed subproblem (3) for all $K \ge 1$, we set $\varepsilon = \frac{C_3}{\sqrt{n}} \le \frac{\lambda_{\text{tgt}}}{8}$ for some generic small constant C_3 .

Computational Theory

 $\|\theta_{\mathcal{S}_{\perp}}^{(t)}\|$

Moreover, we need $\leq T + \log \log \left(3\rho_{s^*+2\tilde{s}}^+ / \varepsilon \right)$ iterations to terminate the prox. Newton update for (4), where $\hat{\theta}^{\{1\}}$ satisfies

$$||\theta_{\mathcal{S}_{\perp}}^{(t)}||_{0} \leq \widetilde{s} \quad \text{and} \quad ||\theta^{(t+1)} - \overline{\theta}^{\{K\}}||_{2} \leq \frac{L_{s^{*}+2\widetilde{s}}}{2\rho_{s^{*}+2\widetilde{s}}^{-}}||\theta^{(t)} - \overline{\theta}^{\{K\}}||_{2}^{2},$$



Statistical Theory

w.h.p.

Theorem 3. Suppose $\{x_i, y_i\}_{i=1}^n$ are generated from GLM satisfying mild conditions with large enough n such that $n \ge C_4 s^* \log d$ and $\beta = C_5/c_{\min}$ is defined in (2), then w.h.p., $\widehat{\theta}^{\{K\}}$ satisfies

- For K-th stage of convex relaxation, we denote $\theta^{\{K\}}$:output solution of (3) $\Rightarrow \omega_{\lambda^{\{K\}}}(\theta^{\{K\}}) \leq \varepsilon$, $\overline{\theta}^{\{K\}}$:uniq. global min. of (3) $\Rightarrow \omega_{\lambda^{\{K\}}}(\overline{\theta}^{\{K\}}) = 0$, $||\overline{\theta}_{S_{\perp}}^{\{K\}}||_{0} \leq \widetilde{s}$.
- **Theorem 1** (Warm Initialization, K = 1). Suppose Assumptions 1 ~ 4 hold. For large enough $T < \infty$, we have for all $t \ge T$,

$$|_{0} \leq \widetilde{s}, \text{ and } ||_{\theta^{(t+1)}} - \overline{\theta}^{\{1\}}||_{2} \leq \frac{L_{s^{*}+2\widetilde{s}}}{2\rho_{s^{*}+2\widetilde{s}}^{-}}||_{\theta^{(t)}} - \overline{\theta}^{\{1\}}||_{2}^{2}.$$

$$||\widehat{\theta}_{\mathcal{S}_{\perp}}^{\{1\}}||_{0} \leq \widetilde{s} \text{ and } ||\widehat{\theta}^{\{1\}} - \theta^{*}||_{2} \leq \frac{18\lambda_{\mathrm{tgt}}\sqrt{s^{*}}}{\rho_{s^{*}+2\widetilde{s}}}.$$

Theorem 2 (Stage K, $K \ge 2$). Suppose Assumptions $1 \sim 4$ hold. Then within each stage $K \geq 2$, for all iterations $t \geq 1$, we have

- Moreover, we need $\leq \log \log \left(3\rho_{s^*+2\tilde{s}}^+/\varepsilon \right)$ iterations to terminate the prox. Newton update for (3), where $\hat{\theta}^{\{K\}}$ satisfies
- $||\widehat{\theta}_{S_{+}}^{\{K\}}||_{0} \leq \widetilde{s} \text{ and } ||\widehat{\theta}^{\{K\}} \theta^{*}||_{2} \leq C_{2}0.7^{K-1}||\widehat{\theta}^{\{1\}} \theta^{*}||_{2} + C_{3}\alpha$
- for some constants C_2 and C_3 , where $\alpha = \|\nabla \mathcal{L}(\theta^*)_{\mathcal{S}}\|_2 +$ $\lambda_{\text{tgt}} \sqrt{\sum_{j \in S} \mathbb{1}(|\theta_j^*| \leq \beta \lambda_{\text{tgt}})} + \varepsilon \sqrt{s^*}.$



- Suppose x_i 's are i.i.d. samples from a *sub-Gaussian* distribution with:
 - bounded eigenvalues of covariance, and
 - other mild conditions.
- Then Assumptions $1\sim 3$ hold





Experiments

Competing algorithms:

Datasets:

Tab. 1 Quantitive timing comparisons for on nonconvex-regularized sparse logistic regression. The average values and standard errors (in parentheses) of timing performance (in seconds) are presented.

	DC+PN	DC+ACD	DC+APG
madelon	$1.51(\pm 0.01)$	$5.83(\pm 0.03)$	$1.60(\pm 0.03)$
	obj value: 0.52	obj value: 0.52	obj value: 0.52
gisette	$5.35(\pm 0.11)$	$18.92(\pm 2.25)$	$207(\pm 2.25)$
	obj value: 0.01	obj value: 0.01	obj value: 0.01
sim_1k	$1.07(\pm 0.02)$	$9.46(\pm 0.09)$	$17.8(\pm 1.23)$
	obj value: 0.01	obj value: 0.01	obj value: 0.01
sim_5k	$4.53(\pm 0.06)$	$16.20(\pm 0.24)$	$111(\pm 1.28)$
	obj value: 0.01	obj value: 0.01	obj value: 0.01
sim_10k	$8.82(\pm 0.04)$	$19.1(\pm 0.56)$	$222(\pm 5.79)$
	obj value: 0.01	obj value: 0.01	obj value: 0.01



Re	eference
[1]	I. Guyon, S. feature selec <i>systems</i> , pag
[2]	T. Zhang. A tion. Journa



• DC+PN: our proposed DC proximal Newton algorithm

• DC+ACD: DC coordinate descent algorithm combined with the active set strategy

• DC+APG: DC accelerated proximal gradient algorithm

• "madelon" (n = 2000, d = 500, [1])• "gisette" (n = 2000, d = 5000, [1])

• simulated datasets: "sim 1k", "sim 5k", and "sim 10k" (n = 1000, d = 1k, 5k, 10k)

Gunn, A. Ben-Hur, and G. Dror. Result analysis of the nips 2003 ection challenge. In Advances in neural information processing ges 545–552, 2005.

Analysis of multi-stage convex relaxation for sparse regularizaal of Machine Learning Research, 11:1081–1107, 2010.

The authors acknowledge support from DARPA YFA N66001-





