

Randomized Block Cubic Newton Method

Nikita Doikov¹

Peter Richtárik^{2, 3, 4}

¹Higher School of Economics, Russia

²King Abdullah University of Science and Technology, Saudi Arabia

³The University of Edinburgh, United Kingdom

⁴Moscow Institute of Physics and Technology, Russia

International Conference on Machine Learning, Stockholm

July 12, 2018

1. Review: Gradient Descent and Cubic Newton methods
2. RBCN: Randomized Block Cubic Newton
3. Application: Empirical Risk Minimization

1. Review: Gradient Descent and Cubic Newton methods
2. RBCN: Randomized Block Cubic Newton
3. Application: Empirical Risk Minimization

Optimization problem:

$$\min_{x \in \mathbb{R}^N} F(x)$$

- ▶ **Main assumption:** gradient of F is Lipschitz-continuous:

$$\|\nabla F(x) - \nabla F(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^N.$$

- ▶ From which we get the **Global upper bound** for the function:

$$F(y) \leq F(x) + \langle \nabla F(x), y - x \rangle + \frac{L}{2}\|y - x\|^2, \quad \forall x, y \in \mathbb{R}^N.$$

- ▶ The Gradient Descent:

$$x^+ = \operatorname{argmin}_{y \in \mathbb{R}^N} \left[F(x) + \langle \nabla F(x), y - x \rangle + \frac{L}{2}\|y - x\|^2 \right] = x - \frac{1}{L}\nabla F(x).$$

Review: Cubic Newton

Optimization problem: $\min_{x \in \mathbb{R}^N} F(x)$.

- ▶ **New assumption:** Hessian of F is Lipschitz-continuous:

$$\|\nabla^2 F(x) - \nabla^2 F(y)\| \leq H\|x - y\|, \quad \forall x, y \in \mathbb{R}^N.$$

- ▶ Corresponding **Global upper bound** for the function:

$$Q(x; y) \equiv F(x) + \langle \nabla F(x), y - x \rangle + \frac{1}{2} \langle \nabla^2 F(x)(y - x), y - x \rangle,$$

$$\text{then } F(y) \leq Q(x; y) + \frac{H}{6} \|y - x\|^3, \quad \forall x, y \in \mathbb{R}^N.$$

- ▶ Newton method with cubic regularization¹:

$$\begin{aligned} x^+ &= \operatorname{argmin}_{y \in \mathbb{R}^N} \left[Q(x; y) + \frac{H}{6} \|y - x\|^3 \right] \\ &= x - \left(\nabla^2 F(x) + \frac{H\|x^+ - x\|}{2} I \right)^{-1} \nabla F(x). \end{aligned}$$

¹Yurii Nesterov and Boris T Polyak. "Cubic regularization of Newton's method and its global performance". In: *Mathematical Programming* 108.1 (2006), pp. 177–205.

Gradient Descent vs. Cubic Newton

Optimization problem: $F^* = \min_{x \in \mathbb{R}^N} F(x)$.

- ▶ $F(x^K) - F^* \leq \varepsilon$, What is K - ?
- ▶ Let F be **convex**: $F(y) \geq F(x) + \langle \nabla F(x), y - x \rangle$.
- ▶ Iteration complexity estimates:
 $K = \mathcal{O}\left(\frac{1}{\varepsilon}\right)$ for GD, and $K = \mathcal{O}\left(\frac{1}{\sqrt{\varepsilon}}\right)$ for CN (much better).
- ▶ But, cost of one iteration: $\mathcal{O}(N)$ for GD and $\mathcal{O}(N^3)$ for CN.

N is huge for modern applications. Even $\mathcal{O}(N)$ is too much!

Recent advances in **block coordinate** methods.

1. Paul Tseng and Sangwoon Yun. “A coordinate gradient descent method for nonsmooth separable minimization”. In: *Mathematical Programming* 117.1-2 (2009), pp. 387–423
2. Peter Richtárik and Martin Takáč. “Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function”. In: *Mathematical Programming* 144.1-2 (2014), pp. 1–38
3. Zheng Qu et al. “SDNA: stochastic dual Newton ascent for empirical risk minimization”. In: *International Conference on Machine Learning*. 2016, pp. 1823–1832

Computationally effective steps, convergence as for full methods.

Aim: To create a second-order method with global complexity guarantees and low cost of every iteration.

1. Review: Gradient Descent and Cubic Newton methods
2. RBCN: Randomized Block Cubic Newton
3. Application: Empirical Risk Minimization

Problem Structure

- ▶ Consider the following decomposition of $F : \mathbb{R}^N \rightarrow \mathbb{R}$:

$$F(x) \equiv \underbrace{\phi(x)}_{\text{twice differentiable}} + \underbrace{g(x)}_{\text{differentiable}}$$

- ▶ For a given space decomposition

$$\mathbb{R}^N \equiv \mathbb{R}^{N_1} \times \dots \times \mathbb{R}^{N_n}, \quad x \equiv (x_{(1)}, \dots, x_{(n)}), \quad x_{(i)} \in \mathbb{R}^{N_i},$$

assume block-separable structure of ϕ :

$$\phi(x) \equiv \sum_{i=1}^n \phi_i(x_{(i)}).$$

- ▶ Block separability for $g : \mathbb{R}^N \rightarrow \mathbb{R}$ is not fixed.

Main Assumptions

Optimization problem: $\min_{x \in Q} F(x)$, where

$$F(x) \equiv \sum_{i=1}^n \phi_i(x_{(i)}) + g(x).$$

- ▶ Every $\phi_i(x_{(i)})$, $i \in \{1, \dots, n\}$ is twice-differentiable and convex, with Lipschitz-continuous Hessian:

$$\|\nabla^2 \phi_i(x) - \nabla^2 \phi_i(y)\| \leq H_i \|x - y\|, \quad \forall x, y \in \mathbb{R}^{N_i}.$$

- ▶ $g(x)$ is differentiable, and for some fixed positive-semidefinite matrices $A \succeq G \succeq 0$ we have bounds, for all $x, y \in \mathbb{R}^N$:
 - ▶ $g(y) \leq g(x) + \langle \nabla g(x), y - x \rangle + \frac{1}{2} \langle A(y - x), y - x \rangle$,
 - ▶ $g(y) \geq g(x) + \langle \nabla g(x), y - x \rangle + \frac{1}{2} \langle G(y - x), y - x \rangle$.
- ▶ $Q \subset \mathbb{R}^N$ is a simple convex set.

Model of the Objective

$$\text{Objective: } F(x) \equiv \sum_{i=1}^n \phi_i(x_{(i)}) + g(x)$$

We want to build a **model** of F .

- ▶ Fix subset of blocks: $S \subset \{1, \dots, n\}$.
- ▶ For $y \in \mathbb{R}^N$ denote by $y_{[S]} \in \mathbb{R}^N$ a vector with zeroed $i \notin S$.
- ▶
$$M_{H,S}(x; y) \equiv F(x) + \langle \nabla \phi(x), y_{[S]} \rangle + \frac{1}{2} \langle \nabla^2 \phi(x) y_{[S]}, y_{[S]} \rangle + \frac{H}{6} \|y_{[S]}\|^3 + \langle \nabla g(x), y_{[S]} \rangle + \frac{1}{2} \langle A y_{[S]}, y_{[S]} \rangle.$$
- ▶ From smoothness: $F(x + y) \leq M_{H,S}(x; y), \forall x, y \in \mathbb{R}^N$
for $H \geq \sum_{i \in S} H_i$.

RBCN: Randomized Block Cubic Newton method

- ▶ Method step: $T_{H,S}(x) \equiv \underset{\substack{y \in \mathbb{R}^N \\ \text{s.t. } x+y \in Q}}{\operatorname{argmin}} M_{H,S}(x; y).$

- ▶ Algorithm:

Initialization: choose $x^0 \in \mathbb{R}^N$, uniform random distribution \hat{S} .

Iterations: $k \geq 0$.

- 1: Sample $S_k \sim \hat{S}$
- 2: Find $H_k > 0$ such that

$$F(x^k + T_{H_k, S_k}(x^k)) \leq M_{H_k, S_k}(x^k; x^k + T_{H_k, S_k}(x^k)).$$

- 3: Make the step: $x^{k+1} \stackrel{\text{def}}{=} x^k + T_{H_k, S_k}(x^k).$

Convergence Results

We want to get: $\mathbb{P}\left(F(x^K) - F^* \leq \varepsilon\right) \geq 1 - \rho$

$\varepsilon > 0$ is required **accuracy level**, $\rho \in (0, 1)$ is **confidence level**.

Theorem 1. General conditions.

$$K = O\left(\frac{1}{\varepsilon} \cdot \frac{n}{\tau} \cdot \left(1 + \log \frac{1}{\rho}\right)\right), \quad \tau \equiv \mathbb{E}[|\hat{S}|].$$

Theorem 2. $\sigma \in [0, 1]$ is a *condition number*. $\sigma \geq \frac{\lambda_{\min}(G)}{\lambda_{\max}(A)} > 0$.

$$K = O\left(\frac{1}{\sqrt{\varepsilon}} \cdot \frac{n}{\tau} \cdot \frac{1}{\sigma} \cdot \left(1 + \log \frac{1}{\rho}\right)\right)$$

Theorem 3. Strongly convex case: $\mu \equiv \lambda_{\min}(G) > 0$.

$$K = O\left(\log\left(\frac{1}{\varepsilon\rho}\right) \cdot \frac{n}{\tau} \cdot \frac{1}{\sigma} \cdot \sqrt{\max\left\{\frac{HD}{\mu}, 1\right\}}\right), \quad D \geq \|x^0 - x^*\|.$$

1. Review: Gradient Descent and Cubic Newton methods
2. RBCN: Randomized Block Cubic Newton
3. Application: Empirical Risk Minimization

ERM problem:

$$\min_{w \in \mathbb{R}^d} \left[P(w) \equiv \sum_{i=1}^n \underbrace{\phi_i(b_i^T w)}_{\text{loss}} + \underbrace{g(w)}_{\text{regularizer}} \right]$$

- ▶ SVM: $\phi_i(a) = \max\{0, 1 - y_i a\}$,
- ▶ Logistic regression: $\phi_i(a) = \log(1 + \exp(-y_i a))$,
- ▶ Regression: $\phi_i(a) = (a - y_i)^2$ or $\phi_i(a) = |a - y_i|$,
- ▶ Support vector regression: $\phi_i(a) = \max\{0, |a - y_i| - \nu\}$,
- ▶ Generalized linear models.

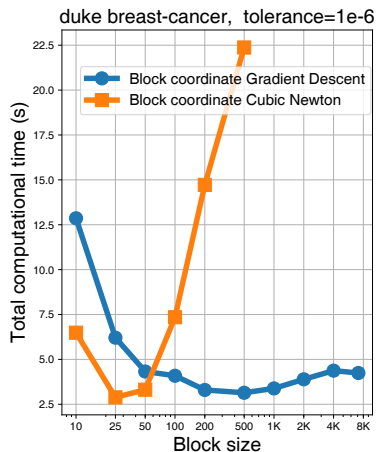
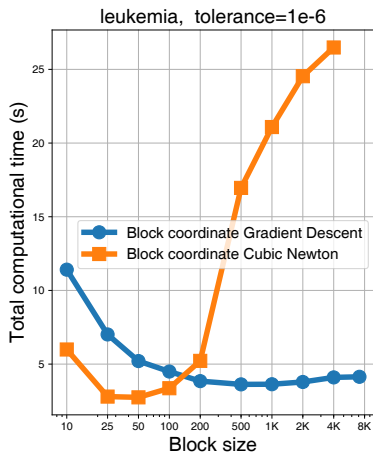
Constrained Problem Reformulation

$$\begin{aligned}\min_{w \in \mathbb{R}^d} P(w) &= \min_{w \in \mathbb{R}^d} \left[\sum_{i=1}^n \phi_i(\underbrace{b_i^T w}_{\equiv \mu_i}) + g(w) \right] \\ &= \min_{\substack{w \in \mathbb{R}^d \\ \mu \in \mathbb{R}^n \\ \underbrace{b_i^T w = \mu_i}_{\equiv Q}}} \left[\underbrace{\sum_{i=1}^n \phi_i(\mu_i)}_{\substack{\text{separable,} \\ \text{twice} \\ \text{differentiable}}} + \underbrace{g(w)}_{\text{differentiable}} \right]\end{aligned}$$

- ▶ Approximate ϕ_i by second-order models with cubic regularization;
- ▶ Treat g as quadratic function;
- ▶ Project onto simple constraints

$$Q \equiv \{w \in \mathbb{R}^d, \mu \in \mathbb{R}^n \mid b_i^T w = \mu_i\}.$$

Proof of Concept: Does second-order information help?



- ▶ Training Logistic Regression, $d = 7129$.
- ▶ Cubic Newton beats Gradient Descent for $10 \leq |S| \leq 50$.
- ▶ Second-order information improves convergence.

Maximization of the Dual Problem

Initial objective: $P(w) \equiv \sum_{i=1}^n \phi_i(b_i^T w) + g(w)$.

We have **Primal** and **Dual** problems:

$$\min_{w \in \mathbb{R}^d} P(w) \geq \max_{\alpha \in \mathbb{R}^n} D(\alpha),$$

introducing Fenchel Conjugate: $f^*(s) \equiv \sup_x [s^T x - f(x)]$, we have

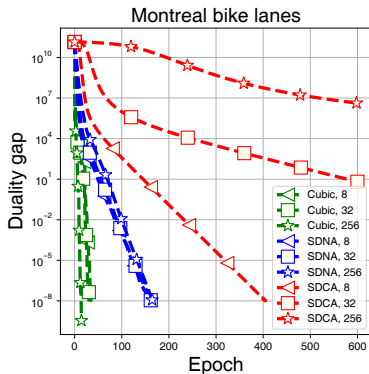
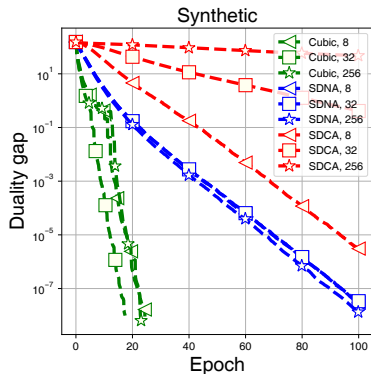
$$D(\alpha) \equiv \underbrace{\sum_{i=1}^n -\phi_i^*(\alpha_i)}_{\text{separable, twice differentiable}} - \underbrace{g^*(-B^T \alpha)}_{\text{differentiable}}.$$

Solve **Dual** problem by our framework:

- ▶ Approximate ϕ_i^* by second-order cubic models;
- ▶ Treat g^* as quadratic function;
- ▶ Project onto $Q \equiv \bigcap_{i=1}^n \text{dom } \phi_i^*$.

Training Poisson Regression

- Solving the dual of Poisson regression.



SDNA: Zheng Qu et al. “SDNA: stochastic dual Newton ascent for empirical risk minimization”. In: *International Conference on Machine Learning*. 2016, pp. 1823–1832

SDCA: Shai Shalev-Shwartz and Tong Zhang. “Stochastic dual coordinate ascent methods for regularized loss minimization”. In: *Journal of Machine Learning Research* 14.Feb (2013), pp. 567–599

Conclusion

New second-order algorithm for convex optimization.

- ▶ Based on cubic regularization.
- ▶ Utilizes problem structure.
- ▶ Does randomized block updates (**computationally cheap**).
- ▶ Has global complexity guarantees.

New Primal-Dual method for Empirical Risk Minimization.

- ▶ Outperforms state-of-the-art in terms of number of data accesses.

Thank you for your attention!

See you at Poster #156.