

Cubic Regularization technical report

Ioannis Tsingalis

Department of Informatics

Aristotle University

September 30, 2021

I. NOTATIONS AND DEFINITIONS

Where needed, the dependence on t is omitted for simplicity. The spectrum of a symmetric $d \times d$ Hessian matrix \mathbf{H} is denoted by $\lambda(\mathbf{H}) = \{\lambda_i(\mathbf{H})\}_{i=1}^d$. Suppose that the eigenvalues are numbered in decreasing order, i.e.,

$$\lambda_1(\mathbf{H}) \geq \lambda_2(\mathbf{H}) \geq \cdots \geq \lambda_d(\mathbf{H}). \quad (1)$$

If \mathbf{H} is indefinite, i.e.,

$$\lambda_d(\mathbf{H}) < 0 \quad \text{and} \quad \lambda_i(\mathbf{H}) > 0, \quad (2)$$

where $i < d$, then $f(\mathbf{x})$ is non-convex.

II. MAIN ALGORITHM

In this section, the Cubic Regularization (CR) technique is described. Let

$$T_M(x) \triangleq \arg \min_y \left[\langle g_x, y - x \rangle + \frac{1}{2} \langle y - x, H_x(y - x) \rangle + \frac{M}{6} \|y - x\|_2^3 \right], \quad (3)$$

where $g_x = \nabla f(x)$ and $H_x = \nabla^2 f(x)$. For a specific time point t we have $g_{x_t} = \nabla_x f(x_t)$ and $H_{x_t} = \nabla^2 f(x_t)$, and $x_{t+1} = T_M(x_t)$.

Algorithm 1 Cubic Regularization of Newton Method

Input: $x_0 \in \mathbb{R}^n$, L_0

- 1: Define $T_M: x \rightarrow T_M(x)$
 - 2: Initialize $t \leftarrow 0$, $M_0 \leftarrow L_0$
 - 3: **do**
 - 4: **while** $f(T_{M_t}(x_t)) > f(x_t)$ **do**
 - 5: $M_t = 2M_t$
 - 6: **end while**
 - 7: $M_{t+1} = M_t$
 - 8: $x_{t+1} = T_{M_t}(x_t)$
 - 9: **while not** $\mu_{M_k}(x_{t+1}) < \epsilon^a$
-

$$\mu_M(x) = \max \left\{ \sqrt{\frac{2}{L+M} \|\nabla_x f(x)\|}, -\frac{2}{2L+M} \lambda_d(\nabla_x^2 f(x)) \right\}$$

III. SOLVING $T_M(x)$

In this section, the procedure to obtain the solution of (3) is described. Subscripts x and t are removed for notational simplicity. Setting $h = y - x$, (3) can be written as

$$\min_{h \in \mathbb{R}^n} \underbrace{\left[g^T h + \frac{1}{2} h^T H h + \frac{M}{6} \|h\|_2^3 \right]}_{m_M(h)}. \quad (4)$$

Even though (4) is non-convex, it is shown in [2, Section 5.1] that the minimizer of (4) is a global minimizer. This follows from the equivalence of

$$\max_{r \in \mathcal{D}} \underbrace{\left[-\frac{1}{2} \langle g, \left(H + \frac{Mr}{2} \right)^{-1} g \rangle - \frac{M}{12} r^3 \right]}_{u(r)}, \quad (5)$$

21 with (4), where $\mathcal{D} = \{r \in \mathbb{R} \mid H + \frac{Mr}{2}I \succ 0, r \geq 0\}$. Assume that r^* is the minimizer of (5). Then the minimizer of (4) is
22 $h^* = h(r^*)$. r^* is obtained by solving $u'(\lambda) = 0$ which is equivalent to

$$\psi(r) = \|h(r)\| - r = 0, \quad (6)$$

23 with

$$h(r) = -\left(H + \frac{Mr}{2}I\right)^{-1}g, \quad r \in \mathcal{D}. \quad (7)$$

24 Setting $\lambda = \frac{M}{2}r$, (6) and (7) can be written as

$$\psi(\lambda) = \|h(\lambda)\| - \frac{2\lambda}{M} = 0 \quad (8)$$

25 and

$$h(\lambda) = -(H + \lambda I)^{-1}g = -H(\lambda)^{-1}g, \quad (9)$$

26 respectively, with $H(\lambda) = H + \lambda I$. Thus, to obtain r^* , the non-linear equation in (8) has to be solved. In [1, Section 7], it is
27 shown that solving (8) directly has drawbacks and the solution of

$$\phi(\lambda) = \frac{1}{\|h(\lambda)\|} - \frac{M}{2\lambda} = 0 \quad (10)$$

28 is preferred.

29 **Lemma III.1.** Suppose $g \neq 0$. Then, the function $\phi(\lambda)$ is strictly increasing, when $\lambda > \lambda_n$ and concave. Its first derivative is

$$\phi'(\lambda) = -\frac{\langle h(\lambda), \partial_\lambda h(\lambda) \rangle}{\|h(\lambda)\|_3^2} + \frac{M}{2\lambda^2}, \quad (11)$$

30 where

$$\partial_\lambda h(\lambda) = -H(\lambda)^{-1}h(\lambda). \quad (12)$$

31 Given the factorization $H(\lambda) = LL^T$, the numerator of $\phi'(\lambda)$ can be written as $\langle h(\lambda), \partial_\lambda h(\lambda) \rangle = -\|w\|^2$, where w is the
32 solution of $Lw = h(\lambda)$. \diamond

33 Proof. A procedure similar to that in [1, Lemma 7.3.1] can be followed to get the results. \square

34 To solve (10), a variety of methods can be applied [1, Section 7]. For example, when the eigensolution is cheap [1, Algorithm
35 7.3.6] can be applied. Algorithm 2 is [1, Algorithm 7.3.6] where $\phi(\lambda)$ and $\phi'(\lambda)$ are defined in (10) and (11), respectively.

Algorithm 2 Subproblem solution of Cubic Regularization of Newton Method

Input: H, g, M
Output: s

```

1: Initialize  $\kappa_{\text{easy}} \in (0, 1)$ 
2: repeat
3:   if  $H \succ 0$  then
4:      $\lambda \leftarrow 0$ 
5:   else
6:      $\lambda \leftarrow -\lambda_d^+$                                  $\triangleright \lambda_d^+$  is barely smaller than  $\lambda_d$ 
7:   end if
8:   Factorize  $H(\lambda) = LL^T$ 
9:   Solve  $LL^T h(\lambda) = -g$ 
10:  if  $\|h(\lambda)\|_2 \leq \frac{2\lambda}{M}$  then
36:    if  $\lambda$  is 0 or  $\|h(\lambda)\|_2 = \frac{2\lambda}{M}$  then
11:      return  $h(\lambda)$ 
12:    else
13:      else
14:        Solve  $(A - \lambda_1)u_1 = 0$  w.r.t.  $u_1$                  $\triangleright$  Get the corresponding eigenvector.
15:        Solve  $\|h(\lambda) + \alpha u_1\|_2 = \frac{2\lambda}{M}$  w.r.t.  $\alpha$ 
16:         $h(\lambda) \leftarrow h(\lambda) + \alpha u_1$ 
17:        return  $h(\lambda)$ 
18:      end if
19:    else
20:      while  $\left|\|h(\lambda)\|_2 - \frac{2\lambda}{M}\right| < \kappa_{\text{easy}} \frac{2\lambda}{M}$  do
21:        Solve  $Lw = h(\lambda)$ 
22:         $\lambda \leftarrow \lambda - \frac{\phi(\lambda)}{\phi'(\lambda)}$            $\triangleright \phi(\lambda)$  and  $\phi'(\lambda)$  are defined in (10) and (11), respectively
23:      end while
24:      return  $h(\lambda)$ 
25:    end if

```

IV. EXPERIMENTS

A. Toy Examples

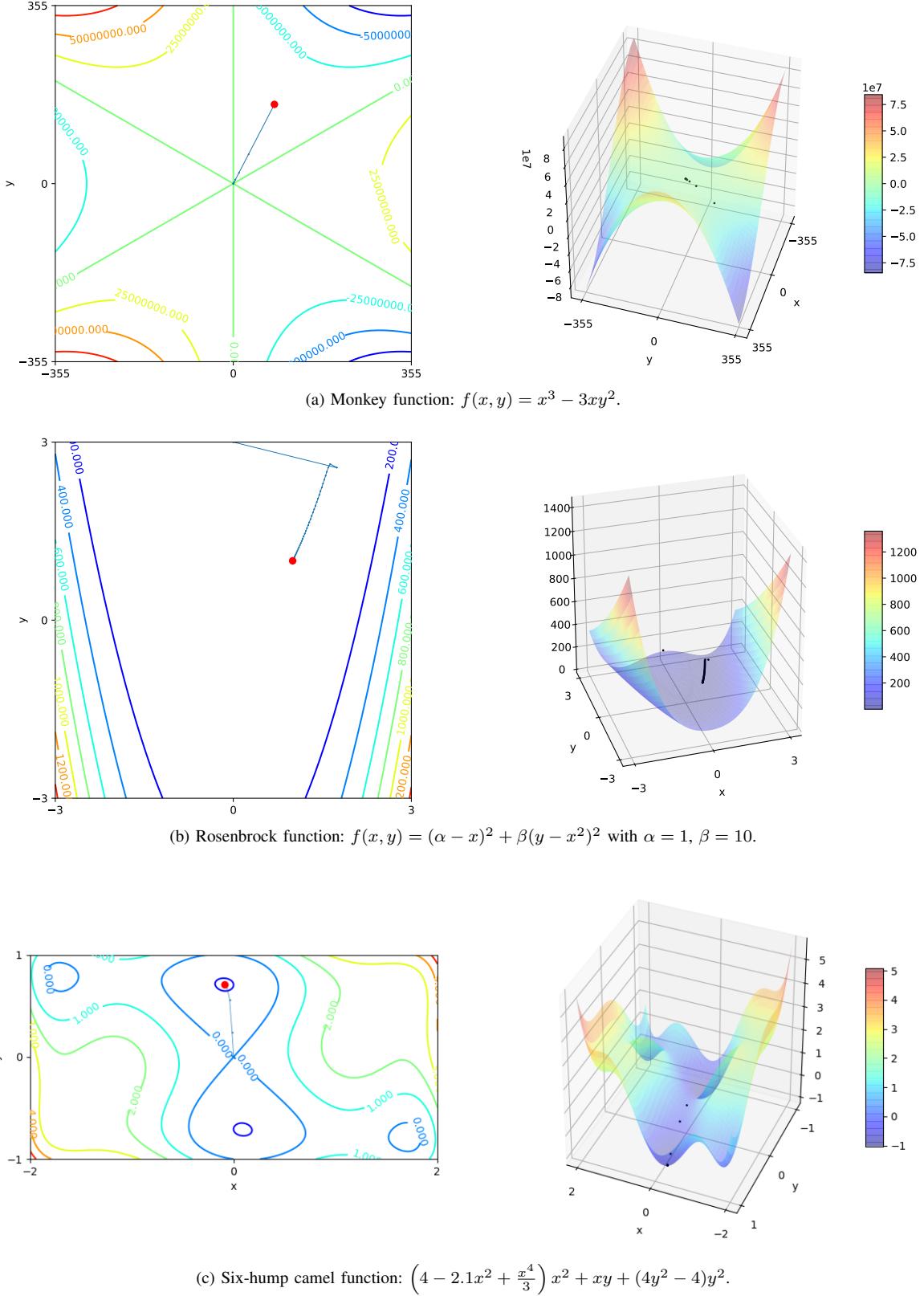


Fig. 1: Application of Cubic Regularization in different functions.

B. Real Data

In this section, CR is applied to binary Logistic Regression with a non-convex regularizer, i.e.,

$$\min_{\mathbf{w} \in \mathbb{R}^d} - \left[\frac{1}{N} \sum_{i=1}^N y_i \log \left(\frac{1}{1 + e^{-w^T x}} \right) + (1 - y_i) \log \left(\frac{e^{w^T x}}{1 + e^{-w^T x}} \right) \right] + \alpha \sum_{j=1}^d \frac{w_j^2}{1 + w_j^2}. \quad (13)$$

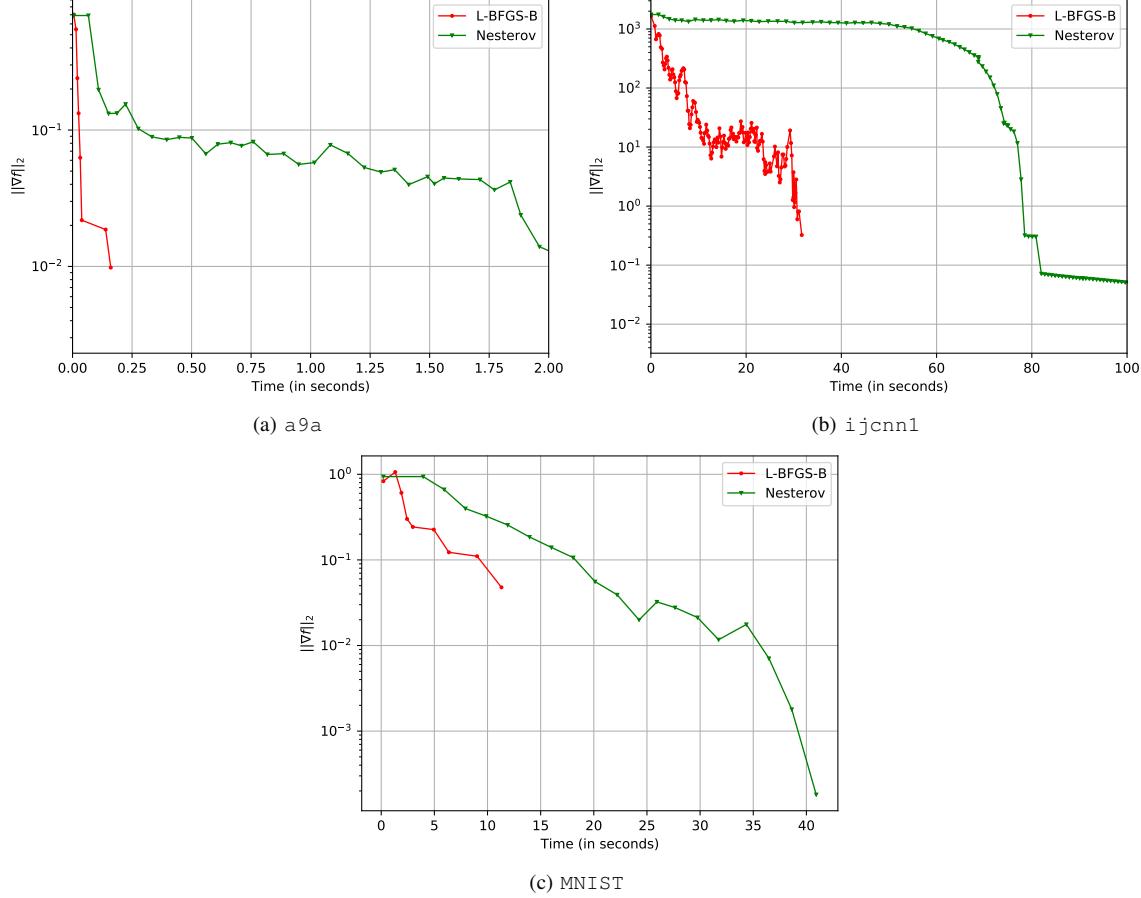


Fig. 2: Average $\|\nabla f\|_2$ for non-convex logistic regression ($\alpha = 0.1$) on a9a ($N=32.561$, $d=123$), covtype ($N= 58.1012$, $d=54$), and MNIST ($N= 60.000$ training / 10.000 testing, $d=784$) after 20 executions.

Methods	Datasets	
	a9a	MNIST
L-BFGS-B	78.71	84.59
CR	77.32	84.73

TABLE I: Average classification accuracy for non-convex logistic regression ($\alpha = 0.1$) on a9a ($N=32.561$, $d=123$), and MNIST ($N= 60.000$ training / 10.000 testing, $d=784$) after 20 executions.

REFERENCES

- [1] A. R Conn, N. IM Gould, and P. L Toint. *Trust region methods*. SIAM, 2000.
 [2] Y. Nesterov and B. T Polyak. Cubic regularization of newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.