

# **Stochastic Cubic Regularization for Fast Nonconvex Optimization**

**Moscow, MIPT, 2020**

**Kyzyl-ool Kezhik,  
Nursultan Kozhogulov  
Ehson Kholzoda**

# Plan

- How does this method differ from classical gradient descent?
- How faster this method than other methods?
- Our experiments and confirmations / refutations

# Gradient descent and Newton method

Classical gradient descent:

$$\mathbf{x}_{t+1}^{\text{GD}} = \operatorname{argmin}_{\mathbf{x}} \left[ f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x} - \mathbf{x}_t) + \frac{\ell}{2} \|\mathbf{x} - \mathbf{x}_t\|^2 \right],$$

Cubic regularized Newton method:

$$\mathbf{x}_{t+1}^{\text{Cubic}} = \operatorname{argmin}_{\mathbf{x}} \left[ f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x} - \mathbf{x}_t) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_t)^\top \nabla^2 f(\mathbf{x}_t) (\mathbf{x} - \mathbf{x}_t) + \frac{\rho}{6} \|\mathbf{x} - \mathbf{x}_t\|^3 \right].$$

# Comparing with other methods

Method	Runtime	Variance Reduction
Stochastic Gradient Descent [Ge et al., 2015]	$\mathcal{O}(\epsilon^{-4}\text{poly}(d))$	not needed
Natasha 2 [Allen-Zhu, 2017]	$\tilde{\mathcal{O}}(\epsilon^{-3.5})^2$	needed
<b>Stochastic Cubic Regularization (this paper)</b>	$\tilde{\mathcal{O}}(\epsilon^{-3.5})$	not needed

# Meta-algorithm

---

**Algorithm 1** Stochastic Cubic Regularization (Meta-algorithm)

---

**Input:** mini-batch sizes  $n_1, n_2$ , initialization  $\mathbf{x}_0$ , number of iterations  $T_{\text{out}}$ , and final tolerance  $\epsilon$ .

```
1: for  $t = 0, \dots, T_{\text{out}}$  do
2:   Sample  $S_1 \leftarrow \{\xi_i\}_{i=1}^{n_1}, S_2 \leftarrow \{\xi_i\}_{i=1}^{n_2}$ .
3:    $\mathbf{g}_t \leftarrow \frac{1}{|S_1|} \sum_{\xi_i \in S_1} \nabla f(\mathbf{x}_t; \xi_i)$ 
4:    $\mathbf{B}_t[\cdot] \leftarrow \frac{1}{|S_2|} \sum_{\xi_i \in S_2} \nabla^2 f(\mathbf{x}_t, \xi_i)(\cdot)$ 
5:    $\Delta, \Delta_m \leftarrow \text{Cubic-Subsolver}(\mathbf{g}_t, \mathbf{B}_t[\cdot], \epsilon)$ 
6:    $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t + \Delta$ 
7:   if  $\Delta_m \geq -\frac{1}{100} \sqrt{\frac{\epsilon^3}{\rho}}$  then
8:      $\Delta \leftarrow \text{Cubic-Finalsolver}(\mathbf{g}_t, \mathbf{B}_t[\cdot], \epsilon)$ 
9:      $\mathbf{x}^* \leftarrow \mathbf{x}_t + \Delta$ 
10:    break
11:  end if
12: end for
```

**Output:**  $\mathbf{x}^*$  if the early termination condition was reached, otherwise the final iterate  $\mathbf{x}_{T_{\text{out}}+1}$ .

---

# Cubic-Subsolver

---

**Algorithm 2** Cubic-Subsolver via Gradient Descent

---

**Input:**  $\mathbf{g}$ ,  $\mathbf{B}[\cdot]$ , tolerance  $\epsilon$ .

1: **if**  $\|\mathbf{g}\| \geq \frac{\ell^2}{\rho}$  **then**

2:    $R_c \leftarrow -\frac{\mathbf{g}^\top \mathbf{B}[\mathbf{g}]}{\rho \|\mathbf{g}\|^2} + \sqrt{\left(\frac{\mathbf{g}^\top \mathbf{B}[\mathbf{g}]}{\rho \|\mathbf{g}\|^2}\right)^2 + \frac{2\|\mathbf{g}\|}{\rho}}$

3:    $\Delta \leftarrow -R_c \frac{\mathbf{g}}{\|\mathbf{g}\|}$

4: **else**

5:    $\Delta \leftarrow 0, \sigma \leftarrow c' \frac{\sqrt{\epsilon \rho}}{\ell}, \eta \leftarrow \frac{1}{20\ell}$

6:    $\tilde{\mathbf{g}} \leftarrow \mathbf{g} + \sigma \zeta$  for  $\zeta \sim \text{Unif}(\mathbb{S}^{d-1})$

7:   **for**  $t = 1, \dots, \mathcal{T}(\epsilon)$  **do**

8:      $\Delta \leftarrow \Delta - \eta(\tilde{\mathbf{g}} + \mathbf{B}[\Delta] + \frac{\rho}{2}\|\Delta\|\Delta)$

9:   **end for**

10: **end if**

11:  $\Delta_m \leftarrow \mathbf{g}^\top \Delta + \frac{1}{2}\Delta^\top \mathbf{B}[\Delta] + \frac{\rho}{6}\|\Delta\|^3$

**Output:**  $\Delta, \Delta_m$

---

# Cubic-Finalsolver

---

**Algorithm 3** Cubic-Finalsolver via Gradient Descent

---

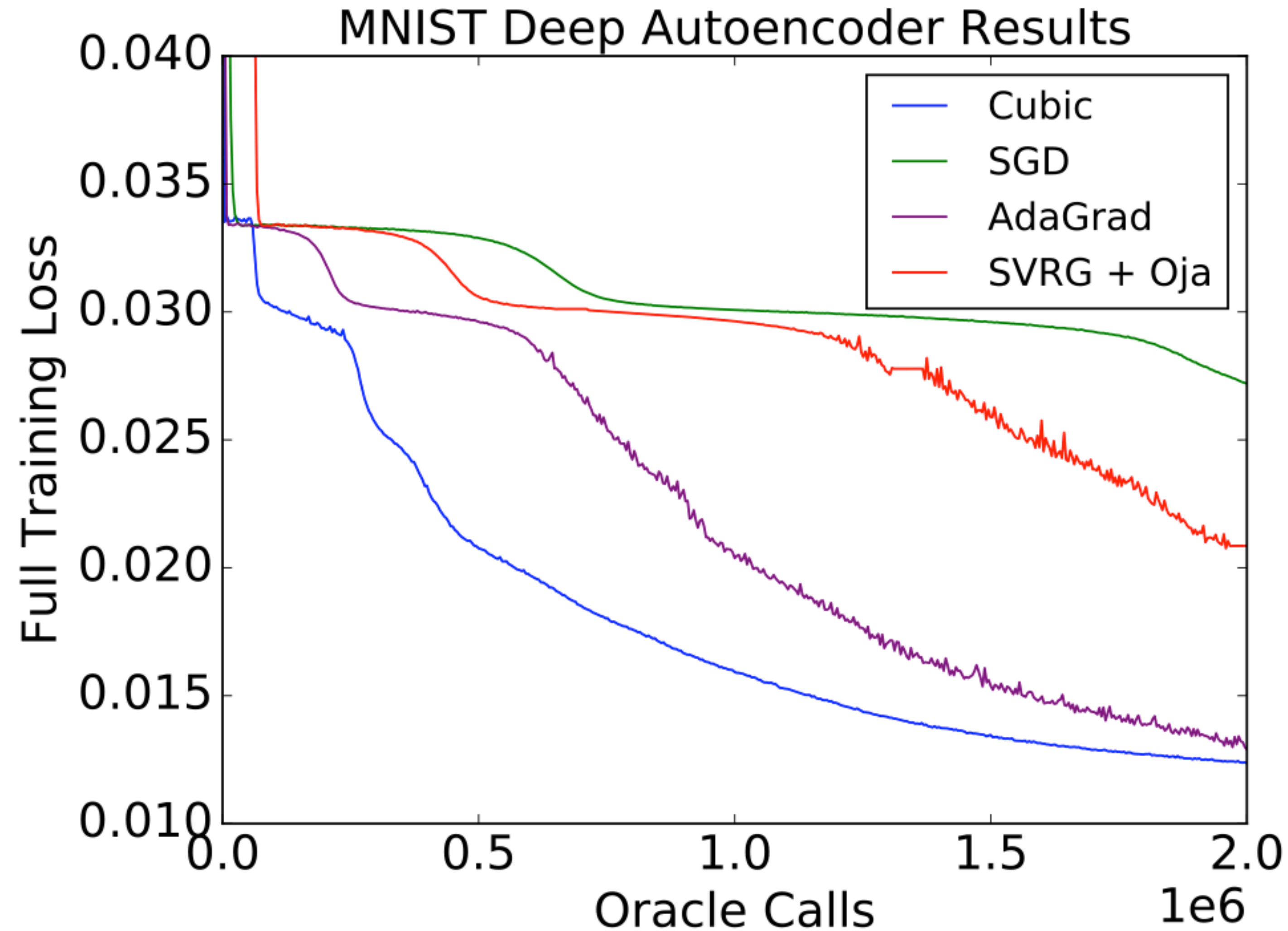
**Input:**  $\mathbf{g}$ ,  $\mathbf{B}[\cdot]$ , tolerance  $\epsilon$ .

```
1:  $\Delta \leftarrow 0$ ,  $\mathbf{g}_m \leftarrow \mathbf{g}$ ,  $\eta \leftarrow \frac{1}{20\ell}$   
2: while  $\|\mathbf{g}_m\| > \frac{\epsilon}{2}$  do  
3:    $\Delta \leftarrow \Delta - \eta \mathbf{g}_m$   
4:    $\mathbf{g}_m \leftarrow \mathbf{g} + \mathbf{B}[\Delta] + \frac{\rho}{2} \|\Delta\| \Delta$   
5: end while
```

**Output:**  $\Delta$

---

# Method compared with other methods





# Our experiments

- To be continued....