

# Incremental Newton Method for Minimizing Big Sums of Functions

Anton Rodomanov

Higher School of Economics  
Bayesian methods research group (<http://bayesgroup.ru>)

28 May 2016  
Seminar on Machine Learning, Voronovo, Russia

# Introduction

## Consider the problem

$$\text{Find } f^* = \min_{x \in \mathbb{R}^d} f(x) \quad \text{with } f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x),$$

## Example (Empirical risk minimization):

- ▶ We are given observations  $a_i$  (and possibly their labels  $\beta_i$ ).
- ▶ Goal: find optimal parameters  $x^*$  of a parametric model.
- ▶ Linear regression ( $a_i \in \mathbb{R}^d$ ,  $\beta_i \in \mathbb{R}$ ):

$$f(x) = \frac{1}{n} \sum_{i=1}^n \|\langle a_i, x \rangle - \beta_i\|^2$$

- ▶ Logistic regression ( $a_i \in \mathbb{R}^d$ ,  $\beta_i \in \{-1, 1\}$ ):

$$f(x) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-\beta_i \langle a_i, x \rangle))$$

- ▶ Neural networks, SVMs, CRFs etc.

# Preliminaries

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x), \quad f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$

## Assumptions:

- ▶ Each function  $f_i$  is *smooth*, i.e. continuously differentiable.
- ▶ Each function  $f_i$  is *convex*:

$$f_i(y) \geq f_i(x) + \langle \nabla f_i(x), y - x \rangle, \quad \forall x, y \in \mathbb{R}^d.$$

- ▶ Function  $f$  is *strongly convex*, i.e. there exists  $\mu > 0$  such that
$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2, \quad \forall x, y \in \mathbb{R}^d.$$

Strong convexity of  $f$  implies existence of a unique  $x^* : f(x^*) = f^*$ .  
We consider iterative methods which produce  $\{x^k\}_{k \geq 0} : x^k \rightarrow x^*$ .

# Gradient descent and big data

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x)$ ,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .

**Gradient descent:**

$$x^{k+1} = x^k - \gamma_k \nabla f(x^k)$$

$$\nabla f(x^k) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k)$$

Here  $\gamma_k \in \mathbb{R}_{++}$  is a step length.

**Note:**

- ▶ Computation of  $\nabla f(x^k)$  requires  $\geq O(nd)$  operations.
- ▶ When  $n$  is very large, this may take a lot of time. Example:  
 $n = 10^8$ ,  $d = 1000 \Rightarrow$  evaluating  $\nabla f(x^k)$  takes  $\approx 2$  minutes.
- ▶ We need a method with less expensive iterations.

# Stochastic gradient descent

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x)$ ,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .

**Stochastic Gradient Descent (SGD):**

Choose  $i_k \in \{1, \dots, n\}$  uniformly at random

$$x^{k+1} = x^k - \gamma_k \nabla f_{i_k}(x^k).$$

Here  $\gamma_k \in \mathbb{R}_{++}$  is a step length.

**Motivation:**  $\mathbb{E}_{i_k} [\nabla f_{i_k}(x^k)] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) = \nabla f(x^k)$ , i.e., on average, SGD makes a step in the right direction.

**Note:**

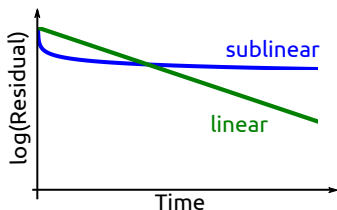
- ▶ Now we only need to compute one gradient instead of  $n$ .
- ▶ Iteration complexity:  $\geq O(d)$ . Independent of  $n$ !

# Gradient descent vs SGD: Which one is better?

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x)$ ,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .

## Iteration cost:

- ▶ Gradient descent:  $O(nd)$ .
- ▶ SGD:  $O(d)$ .



## Convergence rate:

- ▶ Gradient descent: *linear*,  $f(x^k) - f^* = O(c^k)$ ,  $c \in (0, 1)$ .
- ▶ SGD: *sublinear*,  $f(x^k) - f^* = O(\frac{1}{k})$ .

## Discussion:

- ▶ If small accuracy is needed, sublinear is better.
- ▶ If high accuracy is needed, linear is better.
- ▶ Is there a method with  $O(d)$  iteration cost and linear rate?

## Stochastic average gradient [Schmidt et al., 2013]

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x)$ ,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .

### Stochastic Average Gradient (SAG):

Choose  $i_k \in \{1, \dots, n\}$  uniformly at random

Update  $y_i^k = \begin{cases} \nabla f_i(x^k) & \text{if } i = i_k \\ y_i^{k-1} & \text{otherwise} \end{cases}$

$$x^{k+1} = x^k - \gamma_k g^k, \quad \text{where } g^k = \frac{1}{n} \sum_{i=1}^n y_i^k.$$

Here  $\gamma_k \in \mathbb{R}_{++}$  is a step length.

Note that  $g^k = g^{k-1} + \frac{1}{n}(y_{i_k}^k - y_{i_k}^{k-1})$  where  $i = i_k$ .

### Discussion:

- ▶ Variance reduction:  $\mathbb{E}[\|g_k - \nabla f(x^k)\|^2] \rightarrow 0$ .
- ▶ Iteration cost:  $O(d)$  if  $y_i^k$  and  $g_k$  are stored in memory.
- ▶ Convergence rate: linear!

# What if we want more than just a linear rate?

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x), \quad f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$

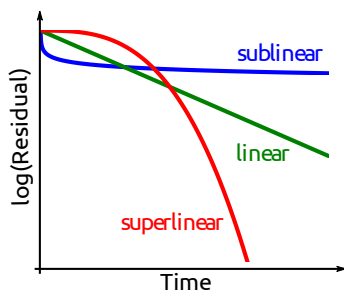
**So far:** methods with sublinear or linear convergence rate.

**Question:** what about superlinear?

**Motivation:** superlinear provides the highest accuracy.

**Convergence rates:**  $r_k := f(x^k) - f^*$

- ▶ Linear:  $r_{k+1} \leq cr_k, \quad c \in (0, 1)$
- ▶ Sublinear:  $r_{k+1} \leq c_k r_k, \quad c_k \uparrow 1$
- ▶ Superlinear:  $r_{k+1} \leq c_k r_k, \quad c_k \downarrow 0.$





# Newton method

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x)$ ,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .

**Newton method:**  $x^{k+1} = x^k - [\nabla^2 f(x^k)]^{-1} \nabla f(x^k)$

$$\nabla f(x^k) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k), \quad \nabla^2 f(x^k) = \frac{1}{n} \sum_{i=1}^n \nabla^2 f_i(x^k).$$

**Interpretation:**

- ▶ Form the second-order Taylor expansion of  $f$  around  $x^k$ :

$$m^k(x) := f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle \nabla^2 f(x^k)(x - x^k), x - x^k \rangle.$$

- ▶ Choose  $x_{k+1}$  as the minimum of the model  $m^k$ :

$$x_{k+1} = \operatorname{argmin}_{x \in \mathbb{R}^d} m^k(x) = x_k - [\nabla^2 f(x^k)]^{-1} \nabla f(x^k).$$

**Convergence rate:** superlinear,  $\|x_{k+1} - x^*\| = O(\|x_k - x^*\|^2)$ .

**Note:** This method is not stochastic, iteration cost depends on  $n$ .

**Question:** Is there a stochastic method with a superlinear rate?

# A quick survey of stochastic optimization methods

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x), \quad f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$

**Consider:** methods whose iteration cost is independent of  $n$ .

**Two groups** of stochastic methods:

- ▶ SGD alike:  $x^{k+1} = x^k - \gamma_k B^k \nabla f_{i_k}(x^k).$ 
  - ▶ SGD, oLBFGS [Schraudolph et al., 2007], AdaGrad [Duchi et al., 2011], SQN [Byrd et al., 2014], Adam [Kingma, 2014] etc.
  - ▶ Convergence rate: sublinear, usually  $O(1/k).$
- ▶ Variance reducing.
  - ▶ IAG [Blatt et al., 2007], SAG [Schmidt et al., 2013], SVRG [Johnson & Zhang, 2013], FINITO [Defazio et al., 2014b], SAGA [Defazio et al., 2014a], MISO [Mairal, 2015] etc.
  - ▶ Convergence rate: linear,  $O(c^k).$

**Note:** no stochastic methods with superlinear convergence.

# Incremental Newton Method – 1 [Rodomanov & Kropotov, 2016]

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x)$ ,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .

**Idea:**

- ▶ Build the second-order Taylor approximation of each  $f_i$ :

$$m_i^k(x) := f_i(v_i^k) + \langle \nabla f_i(v_i^k), x - v_i^k \rangle + \frac{1}{2} \langle \nabla^2 f_i(v_i^k)(x - v_i^k), x - v_i^k \rangle.$$

- ▶ Then  $f$  can be approximated with  $m^k(x) := \frac{1}{n} \sum_{i=1}^n m_i^k(x)$ .
- ▶ Choose the next iterate  $x^{k+1}$  as the minimum of  $m^k$ :

$$x^{k+1} = \operatorname{argmin}_{x \in \mathbb{R}^d} m^k(x).$$

- ▶ Update only one  $v_i^k$  at every iteration to keep the iteration cost independent of  $n$ : choose  $i_k \in \{1, \dots, n\}$  and set

$$v_i^k = \begin{cases} x^k & \text{if } i = i_k, \\ v_i^{k-1} & \text{otherwise.} \end{cases}$$

## Incremental Newton Method – 2 [Rodomanov & Kropotov, 2016]

**Model of the objective:**

$$m^k(x) = \frac{1}{n} \sum_{i=1}^n [f_i(v_i^k) + \langle \nabla f_i(v_i^k), x - v_i^k \rangle + \frac{1}{2} \langle \nabla^2 f_i(v_i^k)(x - v_i^k), x - v_i^k \rangle]$$

**Note:**  $m^k$  is a quadratic,

$$m^k(x) = \frac{1}{2} \langle H^k x, x \rangle + \langle g^k - u^k, x \rangle + \text{const},$$

and determined only by the following three quantities:

$$H^k := \frac{1}{n} \sum_{i=1}^n \nabla^2 f_i(v_i^k), \quad g^k := \frac{1}{n} \sum_{i=1}^n \nabla f_i(v_i^k), \quad u^k := \frac{1}{n} \sum_{i=1}^n \nabla^2 f_i(v_i^k) v_i^k.$$

Since only one component is updated at every iteration,

$$H^k = H^{k-1} + \frac{1}{n} \left[ \nabla^2 f_i(v_i^k) - \nabla^2 f_i(v_i^{k-1}) \right],$$

$$g^k = g^{k-1} + \frac{1}{n} \left[ \nabla f_i(v_i^k) - \nabla f_i(v_i^{k-1}) \right],$$

$$u^k = u^{k-1} + \frac{1}{n} \left[ \nabla^2 f_i(v_i^k) v_i^k - \nabla^2 f_i(v_i^{k-1}) v_i^{k-1} \right].$$

## Incremental Newton Method – 3 [Rodomanov & Kropotov, 2016]

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x), \quad f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$

### Incremental Newton Method (NIM):

Take  $i_k = k \bmod n + 1$

$$\text{Update } v_i^k = \begin{cases} x^k & \text{if } i = i_k \\ v_i^{k-1} & \text{otherwise} \end{cases}$$

$$H^k = H^{k-1} + \frac{1}{n} \left[ \nabla^2 f_{i_k}(v_{i_k}^k) - \nabla^2 f_{i_k}(v_{i_k}^{k-1}) \right],$$

$$g^k = g^{k-1} + \frac{1}{n} \left[ \nabla f_{i_k}(v_{i_k}^k) - \nabla f_{i_k}(v_{i_k}^{k-1}) \right],$$

$$u^k = u^{k-1} + \frac{1}{n} \left[ \nabla^2 f_{i_k}(v_{i_k}^k) v_{i_k}^k - \nabla^2 f_{i_k}(v_{i_k}^{k-1}) v_{i_k}^{k-1} \right]$$

Compute  $x^{k+1} = (H^k)^{-1}(u_k - g_k).$

**Note:** Iteration cost is independent of  $n$  if  $v_i^k$  are kept in memory.

## Superlinear convergence rate of NIM

**Problem:**  $f^* = \min_{x \in \mathbb{R}^d} f(x)$ ,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .

**Theorem:** Suppose the Hessians  $\nabla^2 f_i$  are Lipschitz-continuous:

$$\|\nabla^2 f_i(x) - \nabla^2 f_i(y)\| \leq M \|x - y\|, \quad \forall x, y \in \mathbb{R}^d.$$

Assume  $x^*$  is a minimizer of  $f$  with positive definite Hessian:

$$\nabla^2 f(x^*) = \frac{1}{n} \sum_{i=1}^n \nabla^2 f_i(x^*) \succeq \mu I, \quad \mu > 0,$$

and all the initial points  $x^0, \dots, x^{n-1}$  are close enough to  $x^*$ :

$$\|x^i - x^*\| \leq \frac{\mu}{2M}.$$

Then the sequence of iterates  $\{x^k\}_{k \geq n}$  of NIM converges to  $x^*$  at an R-superlinear rate, i.e. there exists  $\{z_k\}_{k \geq 0}$  such that

$$\|x^k - x^*\| \leq z_k, \quad z_{k+1} \leq \left(1 - \frac{3}{4n}\right)^{2^{\lceil k/n \rceil - 1}} z_k.$$

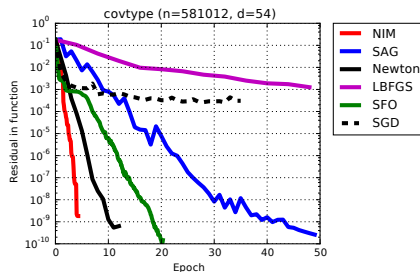
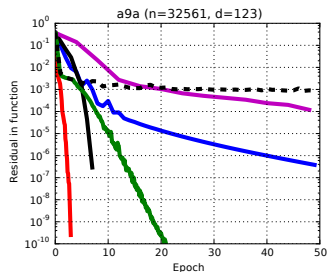
More precisely, the convergence rate is  $n$ -step quadratic:

$$z_{k+n} \leq \frac{M}{\mu} z_k^2.$$

# Evaluation results – 1

**L2-regularized logistic regression:**

$$f(x) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-\beta_i \langle a_i, x \rangle)) + \frac{\mu}{2} \|x\|^2.$$



## Evaluations results – 2

Res	<i>a9a</i> ( $n=32561, d=123$ )					<i>covtype</i> ( $n=581012, d=54$ )				
	NIM	SAG	Newton	LBFGS	SFO	NIM	SAG	Newton	LBFGS	SFO
$10^{-1}$	<b>.01s</b>	.01s	.31s	.05s	.03s	.19s	.33s	.84s	.54s	<b>.04s</b>
$10^{-2}$	<b>.02s</b>	.05s	.56s	.10s	.08s	.51s	.96s	1.78s	1.77s	<b>.25s</b>
$10^{-3}$	.12s	<b>.11s</b>	.73s	.18s	.57s	<b>.72s</b>	1.58s	2.39s	5.67s	1.02s
$10^{-4}$	<b>.15s</b>	.19s	.81s	.43s	.98s	<b>.86s</b>	2.45s	3.09s	10.73s	3.80s
$10^{-5}$	<b>.21s</b>	.36s	.90s	.76s	1.34s	<b>1.20s</b>	3.37s	3.99s	19.07s	5.23s
$10^{-6}$	<b>.24s</b>	.66s	.93s	1.11s	1.57s	<b>1.49s</b>	4.12s	4.57s	31.84s	6.81s
$10^{-7}$	<b>.28s</b>	1.04s	1.00s	1.45s	1.93s	<b>1.69s</b>	4.69s	5.13s	-	8.23s
$10^{-8}$	<b>.31s</b>	1.46s	1.04s	1.82s	2.18s	<b>1.92s</b>	5.90s	6.52s	-	9.86s
$10^{-9}$	<b>.32s</b>	1.90s	1.04s	2.26s	2.46s	<b>2.10s</b>	7.34s	7.64s	-	11.30s
$10^{-10}$	<b>.34s</b>	2.38s	1.04s	2.61s	2.81s	<b>2.12s</b>	9.97s	8.84s	-	12.44s



## Evaluations results – 3

Res	<i>alpha</i> ( $n=500000$ , $d=500$ )				<i>mnist8m</i> ( $n=8100000$ , $d=784$ )			
	NIM	SAG	Newton	LBFGS	NIM	SAG	Newton	LBFGS
$10^{-1}$	1.91s	<b>1.36s</b>	1.6m	4.01s	57.68s	<b>34.91s</b>	47.8m	1.1m
$10^{-2}$	13.37s	<b>6.72s</b>	2.6m	17.68s	<b>1.6m</b>	2.1m	1.4h	5.2m
$10^{-3}$	28.56s	<b>17.73s</b>	3.0m	37.70s	<b>3.2m</b>	3.9m	-	22.9m
$10^{-4}$	36.65s	<b>36.04s</b>	3.4m	58.35s	16.7m	<b>7.1m</b>	-	1.6h
$10^{-5}$	<b>46.66s</b>	1.0m	3.6m	1.4m	<b>26.7m</b>	1.0h	-	-
$10^{-6}$	<b>53.92s</b>	1.5m	4.0m	1.9m	<b>33.5m</b>	-	-	-
$10^{-7}$	<b>57.63s</b>	2.0m	4.0m	2.4m	<b>40.1m</b>	-	-	-
$10^{-8}$	<b>1.0m</b>	2.7m	4.1m	2.8m	<b>46.0m</b>	-	-	-
$10^{-9}$	<b>1.1m</b>	3.5m	4.3m	3.2m	<b>49.6m</b>	-	-	-
$10^{-10}$	<b>1.2m</b>	4.3m	4.7m	3.4m	<b>53.3m</b>	-	-	-

# Conclusions

- ▶ The presented incremental Newton method is the first stochastic method with a superlinear convergence rate.
- ▶ The method can be thought of as a generalization of the classic Newton method to the special case of big sums.
- ▶ It has the same advantages and disadvantages as the classic Newton method:
  - + Fast superlinear convergence rate.
  - Only local convergence is guaranteed.
  - Not applicable to high-dimensional problems.
- ▶ For details, see paper  
A. Rodomanov, D. Kropotov. A Superlinearly-Convergent Proximal Newton-type Method for the Optimization of Finite Sums, ICML 2016.

Thank you!