

# Stochastic Universal Gradient

Evgenii Lagutin

*Optimization Class Project. MIPT*

## Introduction

The universal gradient method is known to be a good approach to numerical optimization when one doesn't have information about the Lipsitz constant of the gradient. This adaptive method adjusts  $L$  at each step of the optimization process and holds the following estimation of the number of calls to the oracle, returning the gradient of the minimized function:

$$N = \inf_{v \in [0,1]} \left( \frac{2L_v R^{1+v}}{\varepsilon} \right)^{\frac{2}{1+v}},$$

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L_v \|y - x\|_2^v, v \in [0,1], L_0 < \infty$$

But this estimation hasn't been transferred on the stochastic case. The purpose of the project is to investigate the effectiveness of the stochastic universal gradient method in practice.

## Algorithm

### Adaptive Stochastic Gradient (Spokoiny's practical variant)

**Input:** lower estimate for the variance of the gradient  $D_0 \leq D$ , accuracy  $0 < \varepsilon < \frac{D_0}{L}$ , starting point  $x_0 \in Q$ , initial guess  $L_{-1} > 0$

- 1: **for**  $k = 0, 1, \dots$  **do**
- 2:   Set  $i_k = 0$ . Set  $r^k = \lceil \frac{2D_0}{L_{k-1}} \varepsilon \rceil$ , generate i.i.d.  $\xi_K^i$ ,  $i = 1, \dots, r^k$
- 3:   **repeat**
- 4:     Set  $L_k = 2^{i_k-1} L_{k-1}$
- 5:     Calculate  $\tilde{g}(x_k) = \frac{1}{r^k} \sum_{i=1}^{r^k} \nabla f(x_k, \xi_k^i)$ .
- 6:     Calculate  $w_k = x_k - \frac{1}{2L_k} \tilde{g}(x_k)$ .
- 7:     Calculate  $\tilde{f}(x_k) = \frac{1}{r^k} \sum_{i=1}^{r^k} f(x_k, \xi_k^i)$  and  $\tilde{f}(w_k) = \frac{1}{r^k} \sum_{i=1}^{r^k} f(w_k, \xi_k^i)$ .
- 8:     Set  $i_k = i_k + 1$ .
- 9:   **until**  $\tilde{f}(w_k) \leq \tilde{f}(x_k) + \langle \tilde{g}(x_k), w_k - x_k \rangle + \frac{2L_k}{2} \|w_k - x_k\|_2^2 + \frac{\varepsilon}{10}$ .
- 10:   Set  $x_{k+1} = w_k$ ,  $k = k + 1$ .
- 11: **end for**

## Optimization of deep neural networks

Let  $g(x)$  be a stochastic gradient of the function being minimized.

In every iteration we have to check if the following inequality is satisfied:

$$f(w) \leq f(x) + \langle g(x), w - x \rangle + \frac{2L}{2} \|w - x\|_2^2 + \frac{\varepsilon}{10}$$

Substituting  $w$  with the its definition expression,  $w = x - \frac{1}{2L} g(x)$

We will get  $f(w) \leq f(x) - \frac{1}{2L} \|g(x)\|_2^2 + \frac{2L}{2} \frac{1}{4L^2} \|g(x)\|_2^2 + \frac{\varepsilon}{10}$

or  $f(w) \leq f(x) - \frac{1}{4L} \|g(x)\|_2^2 + \frac{\varepsilon}{10}$

Consider  $f(x)$  to be a function of a range of matrices and vectors:

$$f(x) = f(W_1, b_1, \dots, W_n, b_n),$$

$$df(W_1, b_1, \dots, W_n, b_n) = \sum_{i=1}^n \left( \frac{\partial f}{\partial W_i} dW_i + \frac{\partial f}{\partial b_i} db_i \right) (W_1, b_1, \dots, W_n, b_n)$$

The goal is to represent  $df$  in this fashion:

$$df(x) = \langle g(x), dx \rangle$$

In this case  $g(x)$  is the gradient.

Let's notice that in case of  $x$  is vector,  $x \in \mathbb{R}^n$ ,  $g(x) \in \mathbb{R}^n$

$$\langle g(x), x \rangle = \sum_{i=1}^n g_i(x) x_i$$

and so we do if  $X$  is a matrix:  $X \in Mat(n \times m)$ ,  $g(X) \in Mat(n \times m)$

$$\langle g(X), X \rangle = \text{tr}(g(X)X) = g(X) \cdot X = \sum_{i=1}^n \sum_{j=1}^m g_{ij}(X) X_{ij}$$

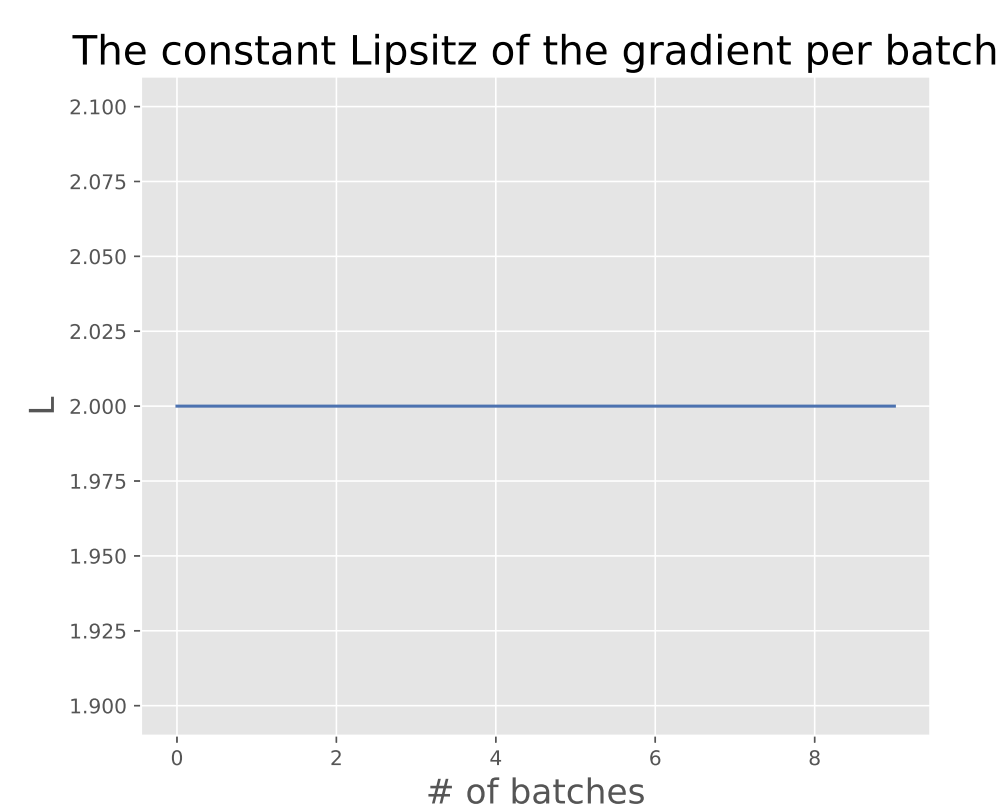
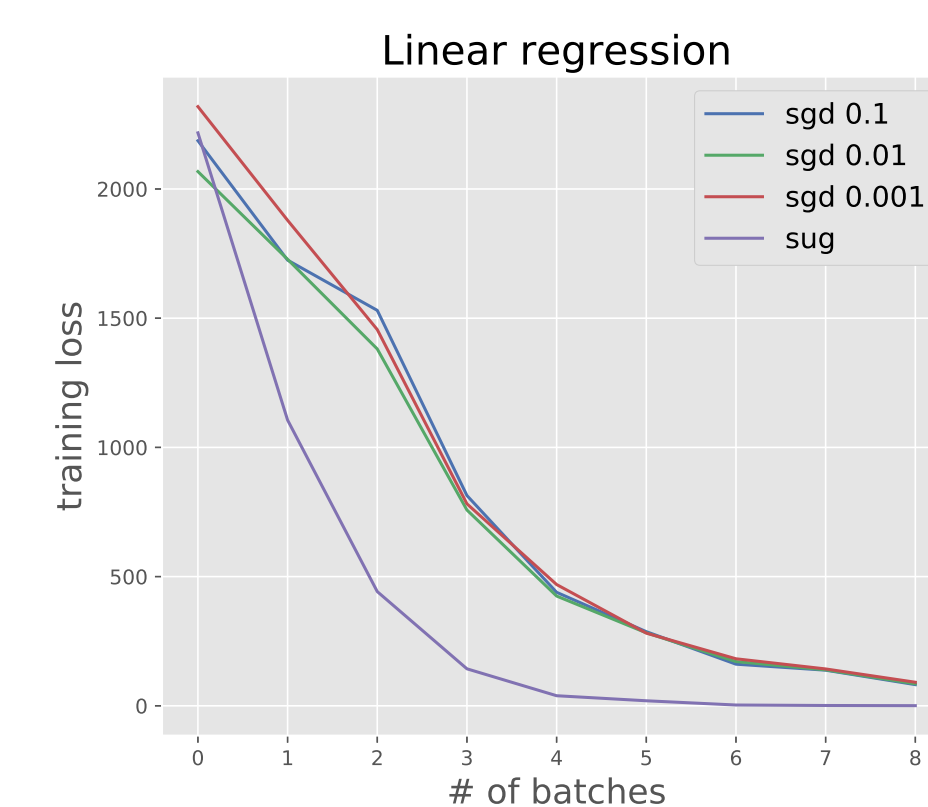
That means we may consider  $X$  as a vector  $(x_{11}, x_{12}, \dots, x_{1m}, x_{21}, \dots, x_{nm})$  of dimension  $nm$  and the result will not change.

Such reasoning allows us to compute the second norm of the gradient in a following way:

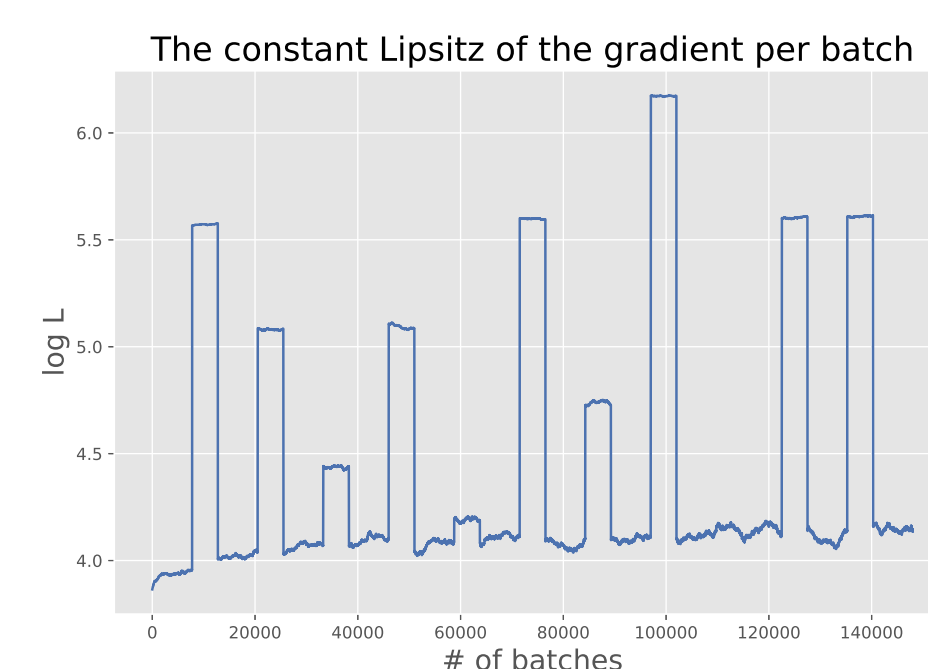
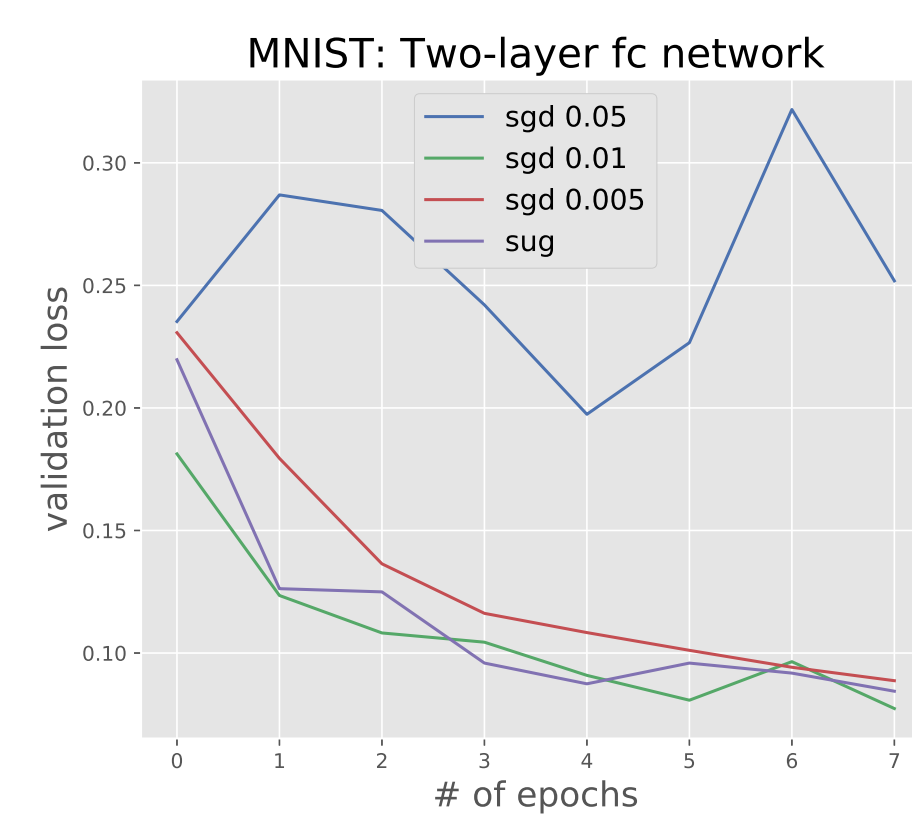
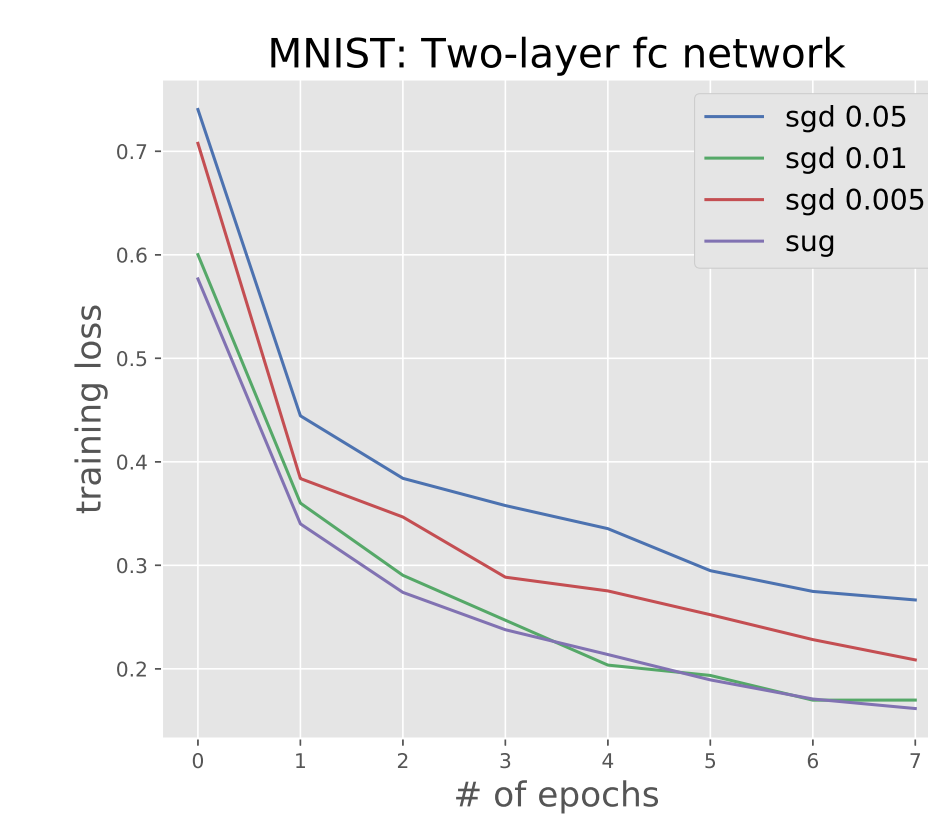
$$\|g(x)\|_2^2 = \|g(W_1, b_1, \dots, W_n, b_n)\|_2^2 = \sum_{i=1}^n (g_{W_i}(x) \cdot g_{W_i}(x) + \langle g_{b_i}(x), g_{b_i}(x) \rangle)$$

## Numerical Experiments

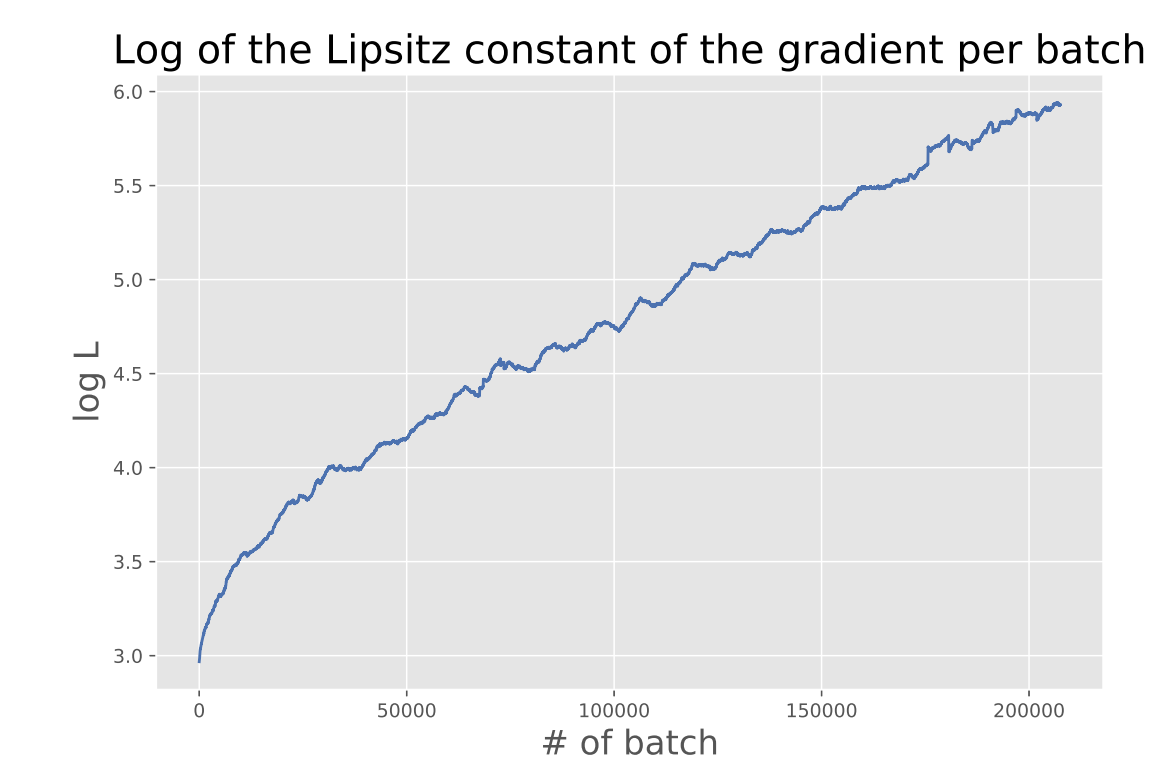
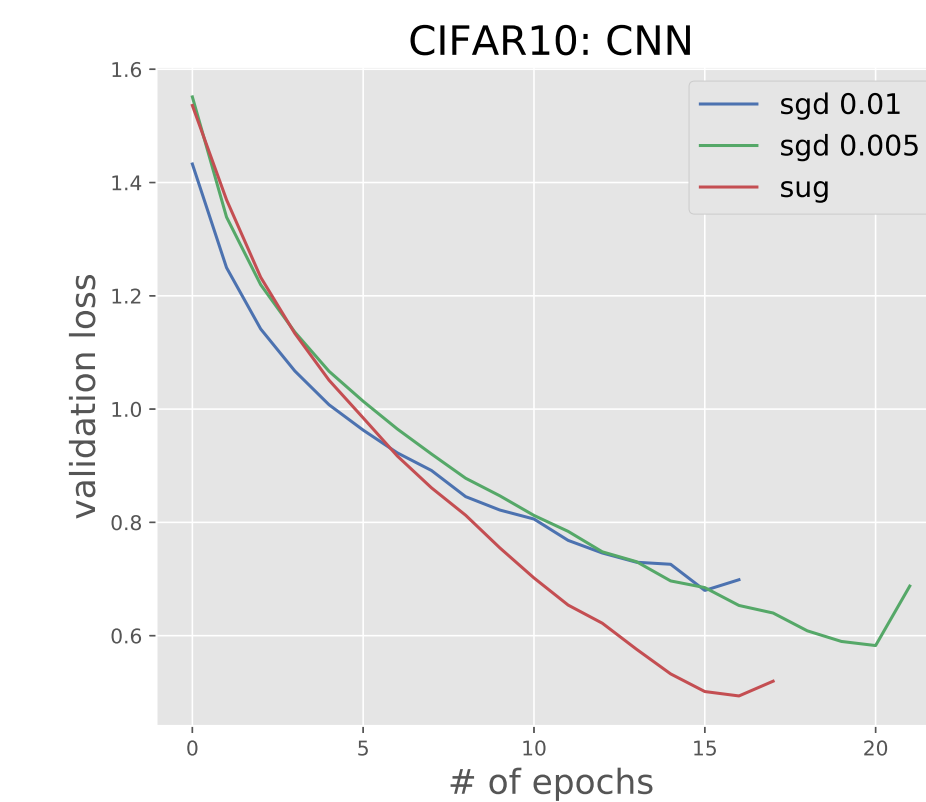
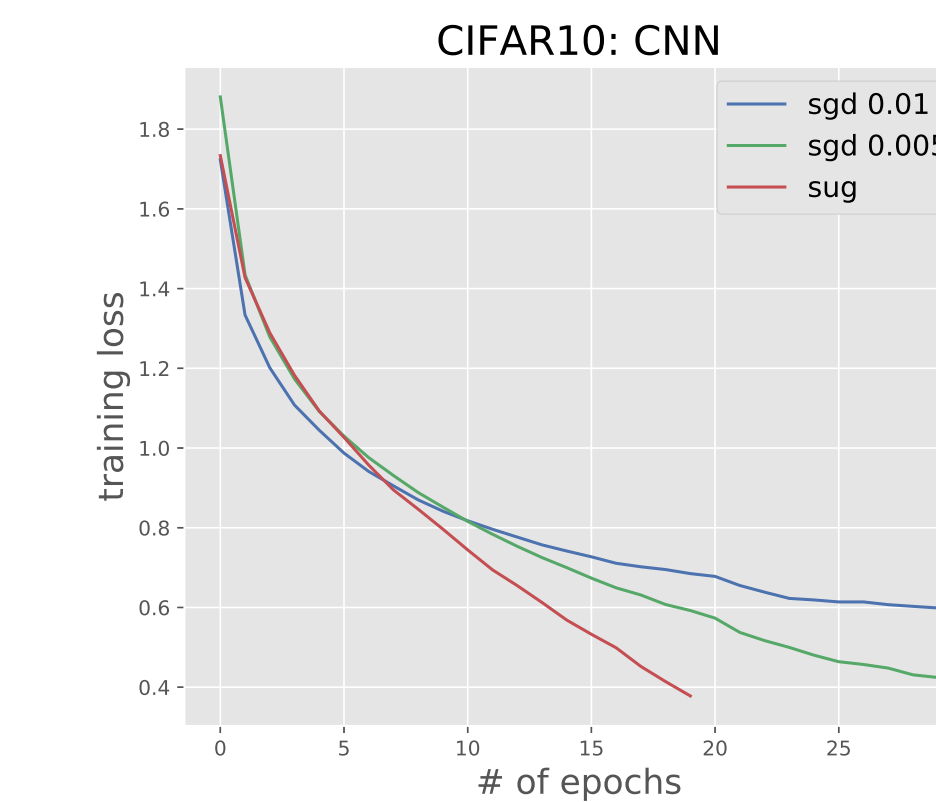
### Linear Regression



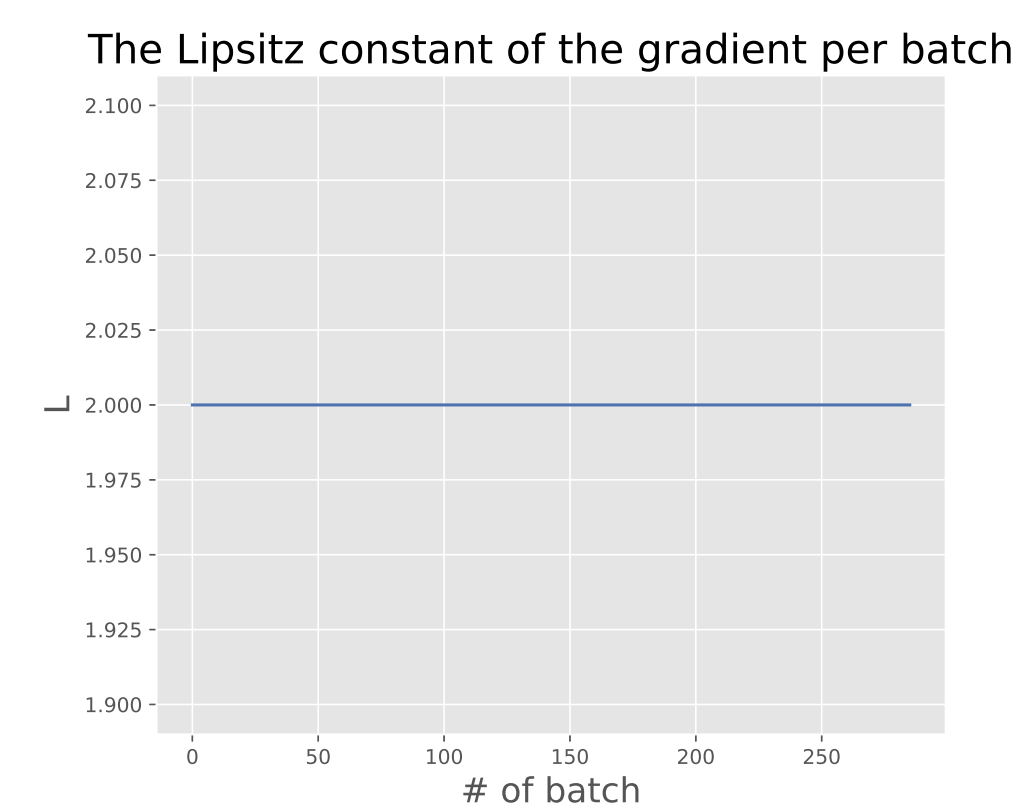
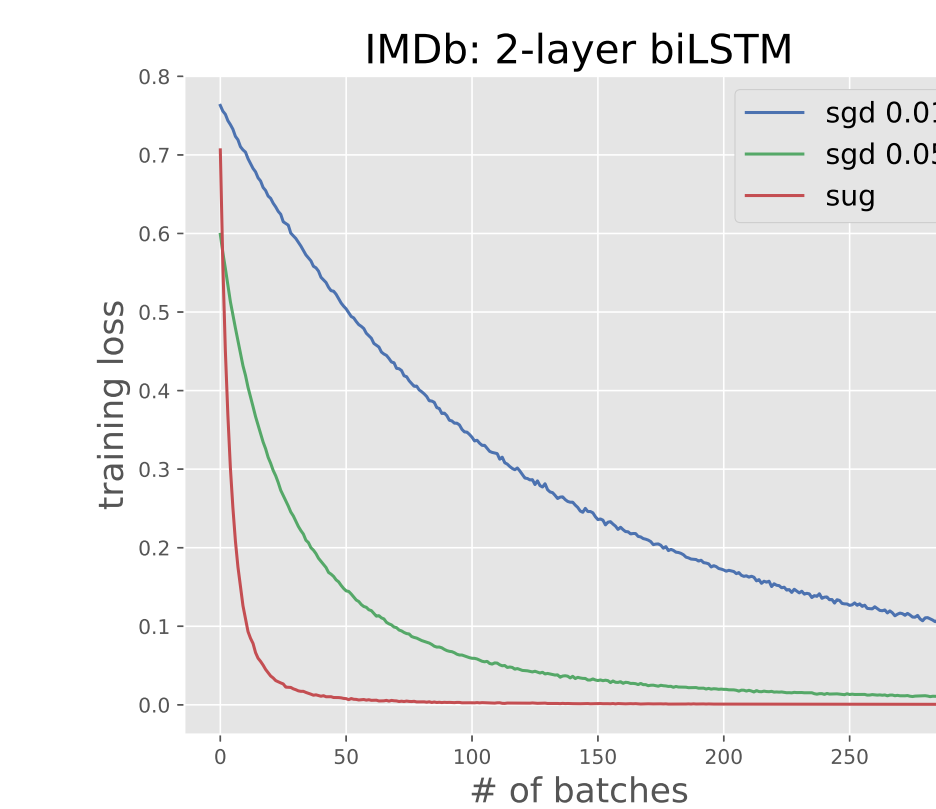
### MNIST



### CIFAR10



### IMDb



## Further actions

Although the method seems to work well especially on complex models, still there are several circumstances that are to be overcome. Doubled time of iteration comparing to the other methods is among them. Besides to maintain the process of optimization stable and avoid undesirable behaviour of the algorithm it was decided to forbid the Lipsitz constant become smaller than 2. That is perhaps an extra constraint not required in theory, so it would be preferable to discard this limit. Still the example with IMDb dataset shows that the method may help significantly decrease the number of optimization steps; the hypothesis is that the solver outperforms others when used to learn complex models. Further it is planned to check this idea on modern neural networks.

## Links

You can watch the project here