

Stochastic Gradient Descent with Momentum and Line Searches

David Nardi

MSc student in Artificial Intelligence, Univeristy of Florence

6th February 2024

Abstract

In recent years, tailored line search approaches have proposed to define the step-size, or learning rate, in SGD-type algorithms for finite-sum problems. In particular, a stochastic extension of standard Armijo line search has been proposed in **bib1**. The development of this kind of techniques is relevant, because it shall allow to enforce a stronger converging behaviour (due to the Armijo condition), similar to that of standard GD, within SGD methods that are commonly employed with large scale training problems.

However, the stochastic line search is not immediately employable when the momentum term is part of the update equation, as the search direction might not be a descent direction (which is a necessary condition for the Armijo condition). This problem is addressed in **bib2**, where a strategy is proposed to guarantee the descent property with momentum.

Contents

1	Introduction	2
1.1	Classification task	2
1.2	Optimization problem	2
2	Mini-batch gradient descent variants	6
2.1	Fixed step-size	6
2.2	Stochastic line search	6
2.3	Fixed momentum term	7
3	Experiments	8
4	Mathematical background	9

1 Introduction

This report summarizes the analysis performed in order to investigate the behaviour of the algorithms retrieved from the scientific literature. The optimization problem that we aim to solve is that of the Logistic Regression with ℓ_2 -regularization term added.

The implemented algorithms are

- Mini-batch Gradient Descent with fixed and decreasing step-size, algorithm 1 on page 6;
- Mini-batch Gradient Descent with Armijo-type line search, algorithm 2 on page 6;
- Mini-batch Gradient Descent with fixed step-size and momentum, algorithm 3 on page 7;
- Mini-batch Gradient Descent with with Armijo-type line search and momentum correction and restart, algorithms 4 and 5 on page 8.

After that, the efficiency of the algorithms is tested on different datasets.

1.1 Classification task

Given the following dataset

$$\mathcal{D} = \{(x^{(i)}, y^{(i)}) \mid x^{(i)} \in \mathcal{X}, y^{(i)} \in \mathcal{Y}, i = 1, 2, \dots, N\}$$

the general Machine Learning (ML) optimization problem in the context of supervised learning is formulated as follows

$$\min_w f(w) = L(w) + \lambda \Omega(w) = \begin{cases} L(w) = \frac{1}{N} \sum_{i=1}^N \ell_i(w) \\ \Omega_{\ell_2} = \frac{1}{2} \|w\|_2^2 \end{cases}$$

where $L(w)$ is called *loss function* and $\Omega(w)$ it's the *regularization term* with its coefficient λ . There are three regularization possible choices, the ℓ_2 regularization was chosen for the problem that we want to address. The vector w contains the model weights associated to the dataset features.

The task performed is the *binary classification*, where $\mathcal{Y} = \{-1, 1\}$ are the allowed values for the response variable, i.e. negative and positive class; the adopted machine learning model is Logistic Regression. Every ML model has its loss function, the logistic regression uses the *log-loss*, for a sample of the dataset the loss function is as follows

$$\ell_i(w) = \log(1 + \exp(-y^{(i)} w^T x^{(i)})) \quad (1)$$

figure 1a on page 5 shows a plot of the loss function where $u = y^{(i)}$ and $v = w^T x^{(i)}$, so the resulting function $\ell(uv) = \log(1 + \exp(-uv))$.

1.2 Optimization problem

Putting together the loss function and the regularization term, we can obtain the optimization problem that we want to solve using Stochastic Gradient Descent (SGD) algorithm variants

$$\min_{w \in \mathbb{R}^{(p+1)}} f(w) = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y^{(i)} w^T x^{(i)})) + \lambda \frac{1}{2} \|w\|^2 \quad (2)$$

where $i = 1, \dots, N$ are the dataset samples, $\mathcal{X} \subseteq \mathbb{R}^{(p+1)}$ where $p + 1$ means that there are p features and the intercept. The $1/N$ term isn't always used, we choose to use that term for scaling issues. We define the matrix associated to the dataset and the model weights as follows

$$X^T = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_p^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_p^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_p^{(N)} \end{pmatrix} \in \mathbb{R}^{N \times (p+1)} \quad x^{(i)} = \begin{pmatrix} 1 \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_p^{(i)} \end{pmatrix} \quad w = \begin{pmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_p \end{pmatrix}$$

the constant column is added for the intercept, also known as *bias*, as the b weight in vector w .

The objective function $f: \mathbb{R}^{(p+1)} \rightarrow \mathbb{R}$ is of class $f \in C^2(\mathbb{R}^{(p+1)})$, we compute the first and second order derivatives

$$\nabla f(w) = \frac{1}{N} X^T r + \lambda w \quad (3)$$

$$\nabla^2 f(w) = \frac{1}{N} X D X^T + \lambda I_{(p+1)} \quad (4)$$

where $r \in \mathbb{R}^N$ is a vector of the same length as the total number of sample, whose elements are $r_i = -y^{(i)} \sigma(-y^{(i)} w^T x^{(i)})$, note that $\sigma(z)$ is the sigmoid function as shown in figure 1b on page 5. $D \in \mathbb{R}^{N \times N}$ is a diagonal matrix whose elements are $d_{ii} = \sigma(y^{(i)} w^T x^{(i)}) \sigma(-y^{(i)} w^T x^{(i)})$ and $I_{(p+1)}$ is the identity matrix with size $p + 1$.

The next proposition allows to address the optimization problem.

Proposition 1. *Problem (2) admits a unique optimal solution.*

Proof. As said, the function is twice continuously differentiable, that allows the use of GD-type algorithms.

(i) Existence of a optimal solution

First thing we prove the existence of an optimal solution, that is the global minimum. The derivatives

where $\sigma(z)$ is the sigmoid function, see figure 1b on page 5 the second order derivative is positive-defined, this because of the square of the euclidean norm, this implies that the objective function is coercive. By a corollary of the Weirstrass theorem we prove that the function admits global minimum.

(ii) Type of optimal solution

□

Remark 1. ...something about convergence speed...

- the hessian matrix is positive defined $\forall w$, this means that the objective function, which is quadratic, is coercive and for the continuity that function admits global minimum, so $f(w)$ has finite inferior limit
- the hessian matrix being positive defined implies also that the objective function is strictly convex (on the other hand the loss function is just convex, due to its hessian matrix being positive semi-defined), this implies that if the global minimum exists, that solution is unique
- a global minimum is a point that satisfy $\nabla f(w^*) = 0$, which is a sufficient condition implied by the convexity of the problem, see figure 1a on page 5

- the ℓ_2 regularization implies that the objective function is strongly convex, this speeds up the convergence
- further more we can assume that $\nabla f(w)$ is Lipschitz-continuous with constant L

$$y = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{pmatrix} \in \{-1, 1\}$$

$$f(w) = \sum_{i=1}^N \log(1 + \exp(-y^{(i)} w^T x^{(i)})) \quad f_i(w) = \log(1 + \exp(-y^{(i)} w^T x^{(i)}))$$

$$\nabla f_i(w) = \frac{1}{N} x^{(i)} r_i + \lambda w$$

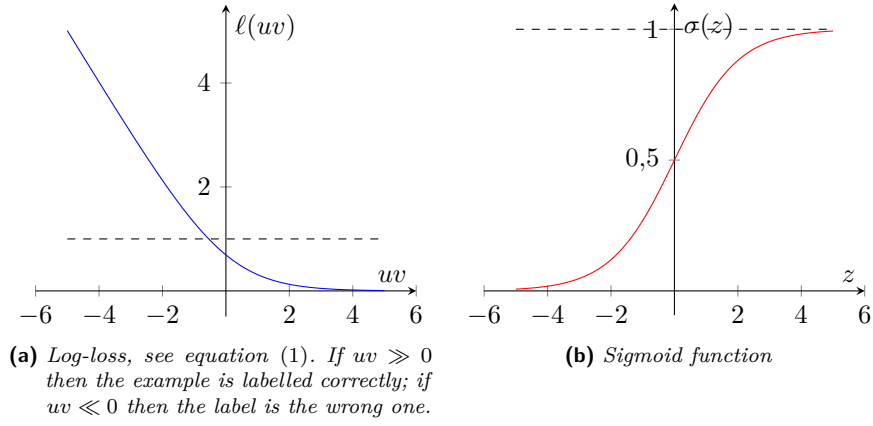


Figure 1: Log-loss and sigmoid function plots

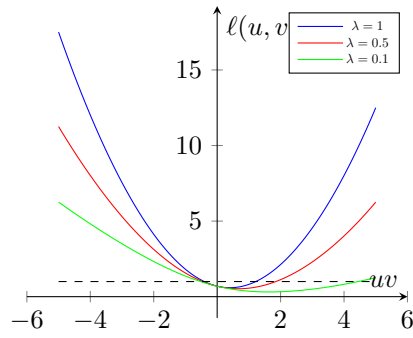


Figure 2

2 Mini-batch gradient descent variants

2.1 Fixed step-size

Algorithm 1: Mini-batch Gradient Descent with fixed or decreasing step-size

```

1 given  $w^0 \in \mathbb{R}^n$ ,  $k = 0$  e  $\{\alpha_k\} \mid \alpha_k = \alpha \vee \alpha_k = \frac{\alpha_0}{k+1}$ 
2 while  $(\|\nabla f(w^k)\| > \varepsilon(1 + |f(w)|))$ 
3   shuffle  $\{1, \dots, N\}$  and split  $B_1, \dots, B_{N/M}$  such that  $1 < |B_t| = M \ll N$ 
4   set  $y_0 = w^k$ 
5   for  $t = 1, \dots, N/M$ 
6     get mini-batch indices from  $B_t$ 
7     approximate true gradient  $\nabla f_{i_t}(w) = \frac{1}{M} \sum_{j \in B_t} \nabla f_j(y_{t-1})$ 
8     compute direction  $d_t = -\nabla f_{i_t}(y_{t-1})$ 
9     make a (internal) step  $y_t = y_{t-1} + \alpha_k d_t$ 
10  end for
11  update weights  $w^{k+1} = y_{N/M}$ 
12  epoch ends  $k = k + 1$ 
13 end while

```

2.2 Stochastic line search

Algorithm 2: Mini-batch Gradient Descent with Armijo line search

```

1 given  $w^0 \in \mathbb{R}^{p+1}$ ,  $\gamma \in (0, 1)$ ,  $\delta \in (0, 1)$ ,  $\alpha_0 \in \mathbb{R}^+$ 
2  $k = 0$ 
3 while  $(\|\nabla f(w^k)\| > \varepsilon(1 + |f(w)|))$ 
4   shuffle  $\{1, \dots, N\}$  and split  $B_1, \dots, B_{N/M}$  such that  $1 < |B_t| = M \ll N$ 
5   set  $y_0 = w^k$ 
6   for  $t = 1, \dots, N/M$ 
7     get mini-batch indices  $i_t$  from  $B_t$ 
8     approximate true gradient  $\nabla f_{i_t}(w) = \frac{1}{M} \sum_{j \in B_t} \nabla f_j(y_{t-1})$ 
9     compute direction  $d_t = -\nabla f_{i_t}(y_{t-1})$ 
10     $\alpha = \text{reset}()$ ,  $q = 0$ 
11    compute potential next step  $y_t = y_{t-1} + \alpha d_t$ 
12    while  $(f_{i_t}(y_t) > f_{i_t}(y_{t-1}) + \gamma \alpha \nabla f_{i_t}(y_{t-1})^T d_t)$ 
13      reduce step-size  $\alpha = \delta \alpha$ 
14      rejections counter  $q = q + 1$ 
15    end while
16    set optimal mini-batch step-size  $\alpha_t = \alpha$ 
17    make a (internal) step  $y_t = y_{t-1} + \alpha_t d_t$ 
18  end for
19  update weights  $w^{k+1} = y_{N/M}$ 
20  epoch ends  $k = k + 1$ 
21 end while

```

2.3 Fixed momentum term

Algorithm 3: Mini-batch Gradient Descent with fixed Momentum term and fixed step-size

```

1 given  $w^0 \in \mathbb{R}^{p+1}$ ,  $\{\alpha_k\} = \alpha$ ,  $\{\beta_k\} = \beta \in (0, 1)$ 
2  $k = 0$ 
3 while  $(\|\nabla f(w^k)\| > \varepsilon(1 + |f(w)|))$ 
4   shuffle  $\{1, \dots, N\}$  and split  $B_1, \dots, B_{N/M}$  such that  $1 < |B_t| = M \ll N$ 
5   set  $y_0 = w^k$ ,  $d_0 = 0$ 
6   for  $t = 1, \dots, N/M$ 
7     get mini-batch indices  $i_t$  from  $B_t$ 
8     approximate true gradient  $\nabla f_{i_t}(w) = \frac{1}{M} \sum_{j \in B_t} \nabla f_j(y_{t-1})$ 
9     compute direction  $d_t = -((1 - \beta)\nabla f_{i_t}(y_{t-1}) + \beta d_{t-1})$ 
10    make a (internal) step  $y_t = y_{t-1} + \alpha_k d_t$ 
11  end for
12  update weights  $w^{k+1} = y_{N/M}$ 
13  epoch ends  $k = k + 1$ 
14 end while

```

Algorithm 4: Mini-batch Gradient Descent with Armijo line search for step-size and Momentum correction

```

1 given  $w^0 \in \mathbb{R}^{p+1}$ ,  $\gamma \in (0, 1)$ ,  $\delta_a \in (0, 1)$ ,  $\alpha_0 \in \mathbb{R}^+$ ,  $\delta_m \in (0, 1)$ ,  $\beta_0 \in (0, 1)$ 
2  $k = 0$ 
3 while  $(\|\nabla f(w^k)\| > \varepsilon(1 + |f(w)|))$ 
4   shuffle  $\{1, \dots, N\}$  and split  $B_1, \dots, B_{N/M}$  such that  $1 < |B_t| = M \ll N$ 
5   set  $y_0 = w^k$ ,  $d_0 = 0$ 
6   for  $t = 1, \dots, N/M$ 
7     get mini-batch indices  $i_t$  from  $B_t$ 
8     approximate true gradient  $\nabla f_{i_t}(w) = \frac{1}{M} \sum_{j \in B_t} \nabla f_j(y_{t-1})$ 
9     compute potential next direction  $d_t = -((1 - \beta)\nabla f_{i_t}(y_{t-1}) + \beta d_{t-1})$ 
10     $q_m = 0$ 
11    while  $(\nabla f_{i_t}(y_{t-1})^T d_t \geq 0)$ 
12      reduce momentum term  $\beta = \delta_m \beta$ 
13      rejections counter  $q_m = q_m + 1$ 
14    end while
15    set optimal mini-batch momentum term  $\beta_t = \beta$ 
16     $\alpha = \text{reset}()$ ,  $q_a = 0$ 
17    compute potential next step  $y_t = y_{t-1} + \alpha d_{t-1}$ 
18    while  $(f_{i_t}(y_t) > f_{i_t}(y_{t-1}) + \gamma \alpha \nabla f_{i_t}(y_{t-1})^T d_t)$ 
19      reduce step-size  $\alpha = \delta_a \alpha$ 
20      rejections counter  $q_a = q_a + 1$ 
21    end while
22    set optimal mini-batch step-size  $\alpha_t = \alpha$ 
23    make a (internal) step  $y_t = y_{t-1} + \alpha_t d_t$ 
24  end for
25  update weights  $w^{k+1} = y_{N/M}$ 
26  epoch ends  $k = k + 1$ 
27 end while

```

Algorithm 5: Mini-batch Gradient Descent with Armijo line search for step-size and Momentum restart

```

1 given  $w^0 \in \mathbb{R}^{p+1}$ ,  $\gamma \in (0, 1)$ ,  $\delta_a \in (0, 1)$ ,  $\alpha_0 \in \mathbb{R}^+$ ,  $\delta_m \in (0, 1)$ ,  $\{\beta_k\} = \beta \in (0, 1)$ 
2  $k = 0$ 
3 while  $(\|\nabla f(w^k)\| > \varepsilon(1 + |f(w)|))$ 
4   shuffle  $\{1, \dots, N\}$  and split  $B_1, \dots, B_{N/M}$  such that  $1 < |B_t| = M \ll N$ 
5   set  $y_0 = w^k$ ,  $d_0 = 0$ 
6   for  $t = 1, \dots, N/M$ 
7     get mini-batch indices  $i_t$  from  $B_t$ 
8     approximate true gradient  $\nabla f_{i_t}(w) = \frac{1}{M} \sum_{j \in B_t} \nabla f_j(y_{t-1})$ 
9     compute potential next direction  $d_t = -((1 - \beta)\nabla f_{i_t}(y_{t-1}) + \beta d_{t-1})$ 
10    if  $(\nabla f_{i_t}(y_{t-1})^T d_t \geq 0)$ 
11      restart direction  $d_t = d_0$ 
12    end if
13     $\alpha = \text{reset}()$ ,  $q_a = 0$ 
14    compute potential next step  $y_t = y_{t-1} + \alpha d_{t-1}$ 
15    while  $(f_{i_t}(y_t) > f_{i_t}(y_{t-1}) + \gamma \alpha \nabla f_{i_t}(y_{t-1})^T d_t)$ 
16      reduce step-size  $\alpha = \delta_a \alpha$ 
17      rejections counter  $q_a = q_a + 1$ 
18    end while
19    set optimal mini-batch step-size  $\alpha_t = \alpha$ 
20    make a (internal) step  $y_t = y_{t-1} + \alpha_t d_t$ 
21  end for
22  update weights  $w^{k+1} = y_{N/M}$ 
23  epoch ends  $k = k + 1$ 
24 end while

```

3 Experiments

4 Mathematical background

Theorem 1 (Weirstrass theorem). *Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous function and $S \subseteq \mathbb{R}^n$ a compact set. Then function f admits global minimum in S .*

Corollary 2 (Sufficient condition). *If function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous and coercive function, then f admits global minimum in \mathbb{R}^n .*

Proposition 2 (Coercivity of a quadratic function). *A quadratic function $f(x) = \frac{1}{2}x^T Qx - c^T x$ is said to be coercive if and only if the symmetric matrix $Q \in \mathbb{R}^{n \times n}$ is positive-defined.*

Proposition 3 (Unique global minimum). *Let $S \subseteq \mathbb{R}^n$ be a convex set, let $f: S \rightarrow \mathbb{R}$ be a strictly convex function. Then the global minimum, if exists, is unique.*

Definition 1 (Convex function). Let $S \subseteq \mathbb{R}^n$ be a convex set, a function $f: S \rightarrow \mathbb{R}$ is said to be convex if the hessian matrix is semi-positive-defined. If the hessian matrix is positive-defined, then the function is strictly convex.