

---

# Stochastic Gradient Descent Report

---

<b>Atou koffi KOUGBANHOUN</b> AIMS AMMI Mbour, Sénégal akougbanhoun@aimsammi.org	<b>Millicent Omondi</b> AIMS AMMI Mbour, Sénégal momondi@aimsammi.org
---	--

**Aime Ouraga**  
AIMS AMMI  
Mbour, Sénégal  
oballou@aimsammi.org

## 1 Introduction

Stochastic Gradient Descent (SGD) is an algorithm that minimizes the loss function by computing its gradient after each training example as opposed to using the whole dataset as seen in gradient descent. This approach makes SGD computationally efficient, especially when using a large data set. Given that SGD chooses one training example from the dataset randomly for each iteration, it suffers from noise. Because of this, there are various modifications to SGD like incorporating momentum, adjusting step sizes, and using different sampling strategies. These help mitigate noisy updates and slow convergence thereby improving performance.

### 1.1 Problem Statement

The most important step in building machine learning models is to solve the optimization problem as defined by the loss function. While standard gradient descent offers reliable method for finding optimal parameters, it may be computationally expensive.

### 1.2 Objectives

There are different variants of SGD. The project's objective is look at the following:

- SGD with constant stepsizes.
- SGD with shrinking stepsizes
- SGD with sampling with/without replacement
- SGD with averaging
- SGD with momentum
- comparing the above with gradient descent.

## 2 Approach

In this part we will look at the loss functions, gradients, stepsizes and the different variants of SGD that we used. We will compare the performance of these SGD variants against each other and with the traditional gradient descent.

## 2.1 Loss functions, gradients and step-sizes

We implemented linear and logistic regression models. These models fall under supervised learning techniques. Linear regression is used for solving regression problems whereas logistic regression is used for classification problems. Our focus will be applying SGD on the linear regression model.

### Loss Functions:

We want to minimize

$$\frac{1}{n} \sum_{i=1}^n \ell(x_i^\top w, b_i) + \frac{\lambda}{2} \|w\|_2^2 \quad (1)$$

where,

$$\ell(z, b) = \frac{1}{2}(b - z)^2 \text{ (least-squares regression)}$$

$$\ell(z, b) = \log(1 + \exp(-bz)) \text{ (logistic regression).}$$

35

We write it as a minimization problem of the form:

$$\frac{1}{n} \sum_{i=1}^n f_i(w)$$

36 where,

$$f_i(w) = \ell(x_i^\top w, y_i) + \frac{\lambda}{2} \|w\|_2^2. \quad (2)$$

## 2.2 Comparison between SGD with constant step sizes and SGD with shrinking step sizes

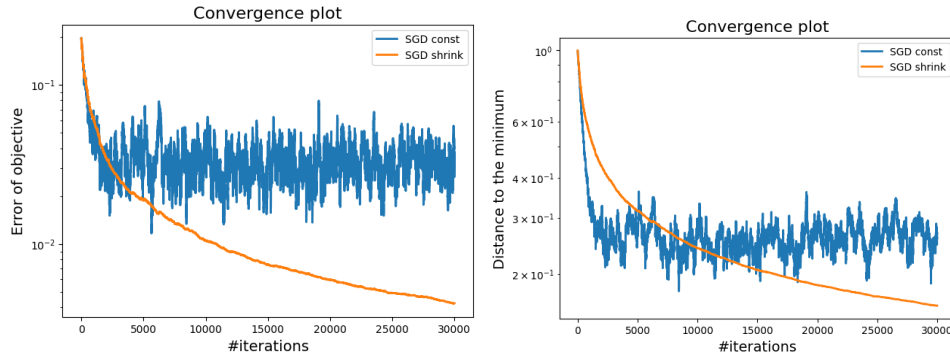


Figure 1: shrink vs constant

From figure 1, we observe that in the initial stages, SGD with constant step sizes shows faster convergence. This is seen by the rapid decrease in the error of the objective and distance to the minimum. However, it oscillates around the minimum and does not converge smoothly.

For the SGD with shrinking step sizes, we see that it starts slower and then gradually reduces the error steadily. We see that it converges toward the minimum steadily without much oscillation. This shows that it has a more stable and consistent reduction in both the error and the distance to the minimum.

From figure 1, we may conclude that SGD with shrinking step sizes reaches the "best" solution as it has a lower error and closer distance to the minimum. On the other hand, SGD with constant step size finds it hard to reach the optimal solution due to the fixed step size thus having continuous oscillations.

**Sampling without replacement:** In this case, each data point is used only once per iteration. It leads to faster convergence as it reduces the variance of the gradient leading to stability and faster convergence within a single data pass.

## 2.3 Comparison between SGD with a switching stepsize and SGD with shrinking step sizes

From figure 2, we observe that SGD with switch starts with larger steps and switches to smaller ones after some iterations.

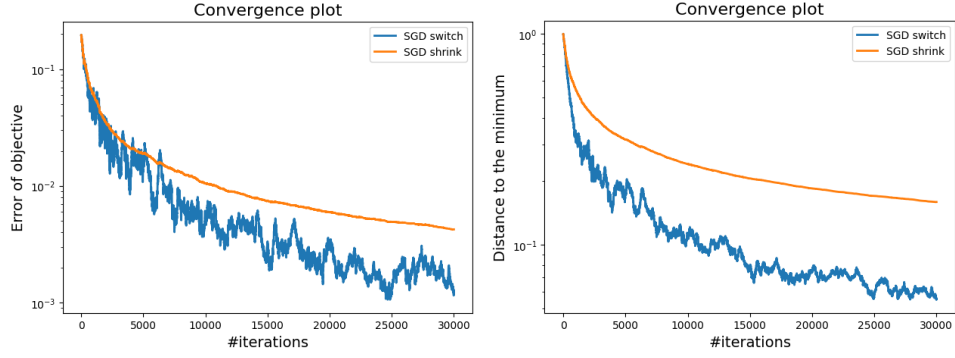


Figure 2: Switch vs shrinking

53 The SGD with shrinking steps also starts with larger steps and then gradually reduces the step sizes.  
 54 As compared to SGD with switch, we observe that the convergence is more gradual and the final  
 55 error is lower.

56 Looking at the two methods, we can deduce that SGD with shrink converges faster at the beginning  
 57 but reaches a higher final error as compared to SGD with shrink. SGD with shrink appears to be more  
 58 stable as it has a smoother convergence curve with fewer fluctuations.

## 59 2.4 Comparison between SGD with switching stepsizes and SGD with averaging

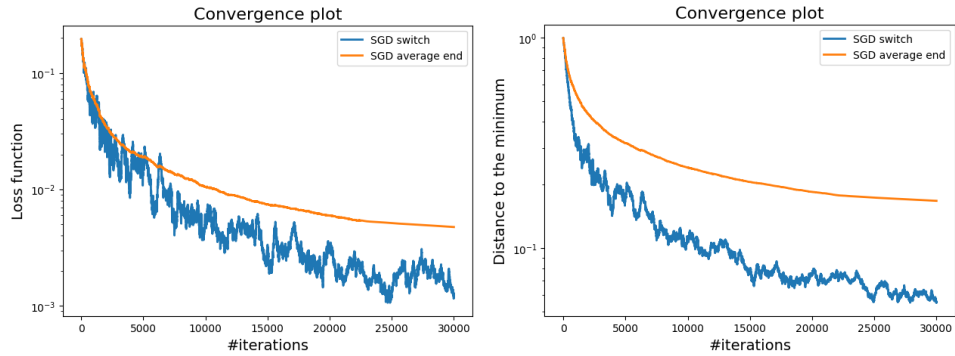


Figure 3: Average vs switch

60 From figure 3, we observe that SGD with switching stepsizes approach exhibits faster initial conver-  
 61 gence at the beginning as large steps allow the algorithm to progress rapidly towards the optimum  
 62 solution. After the steps switch, the convergence rate slows down therefore making the final solution  
 63 to have a higher error as compared to SGD with averaged stepsize.

64 Looking at SGD with averaging, we observe it has a more gradual convergence and has a lower error  
 65 compared to SGD with switching stepsizes.

66 **Averaging only the last  $n$  iterates:** When we average the last  $n$  iterates, we smooth out the noise  
 67 and obtain a more stable and accurate final solution.

68 **Averaging is useful** when we want to improve the final solution in cases where we are dealing with  
 69 noisy gradients near the optimum.

## 70 2.5 Comparison between SGD with switching stepsizes , SGD with momentum and gradient 71 descent.

72 In figure 4, the SGD with switch starts with faster convergence and then reduces. In the end, it's final  
 73 solution has a higher error as compared to the other two methods.

74 As for SGD with momentum, we see that it experiences a sudden drastic reduction in the loss  
 75 function after some iterations. This is depicted by the vertical green line. This is associated with  
 76 the momentum term in its update rule. The momentum term accumulates thereby amplifying the

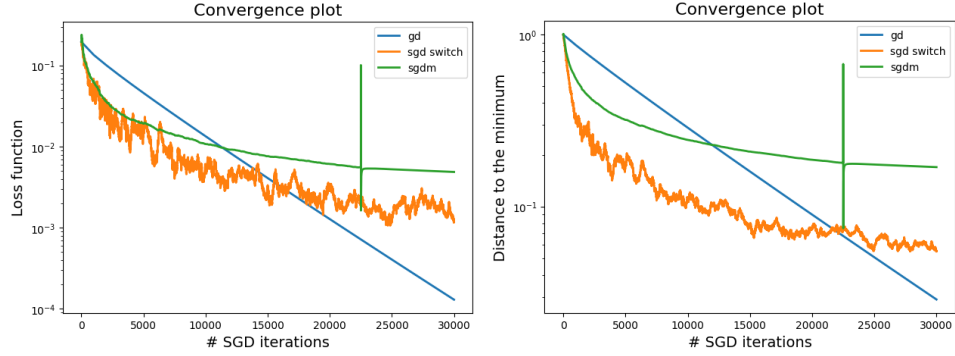


Figure 4: SGD with momentum, decreasing stepsizes and GD

77 gradients allowing the algorithm to escape shallow local minima or saddle points effectively.  
 78 We see that that the gradient descent converges to a lower loss function value as compared to the other  
 79 two methods. This could be due to the fact that the full gradient was computed over the entire dataset  
 80 at each iteration. This implies that the gradient information used in the update is more accurate and  
 81 representative of the true gradient of the objective function. This therefore allows gradient descent to  
 82 converge to a better optimum.

83 **For what momentum parameters does SGDm work well?** momentum value between 0.5 – 0.9. It  
 84 helps the algorithm to adapt quickly to the current gradient direction without being influenced by past  
 85 gradients. In our case we used momentum value of 0.9

86 **Combining all the tricks/variants seen so far, what is the best variant of SGD for this problem?** In  
 87 our case, it would be a combination of decreasing stepsize and momentum. The algorithm will perform  
 88 better as decreasing stepsize will help stabilize the solution near optimum, whereas momentum will  
 89 accelerate convergence.

## 90 2.6 SGD with replacement vs without replacement

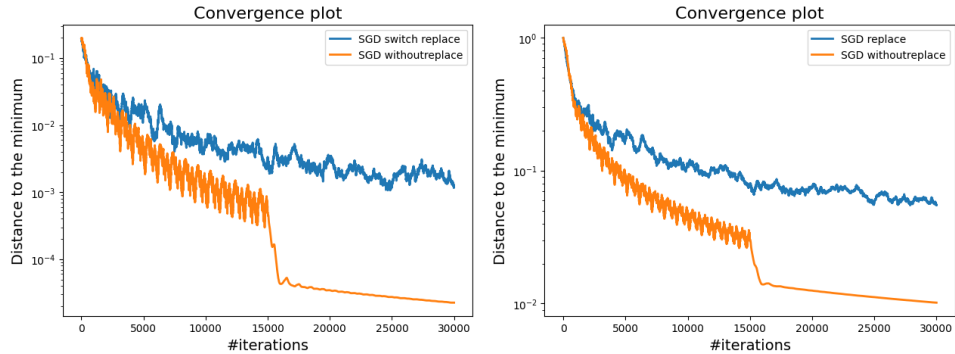


Figure 5: SGD with/without replacement

91 We can see from figure 5, that SGD with replacement converges faster initially, but in the end SGD  
 92 without replacement achieves a better final solution. SGD without replacement shows a more stable  
 93 and smooth convergence trajectory as compared to SGD with replacement, that is , it has fewer  
 94 oscillations.

## 95 2.7 SGD with replacement vs averaging end

96 From figure 6, we observe that SGD without replacement exhibits the quickest initial convergence,  
 97 yet it results in the highest final loss value. In contrast, SGD with averaging begins at a slower pace  
 98 but ultimately outperforms SGD without replacement, leading to a lower final loss. Gradient Descent  
 99 (GD) demonstrates the slowest initial convergence; however, it ultimately reaches the lowest final  
 100 loss value when compared to the other two methods.

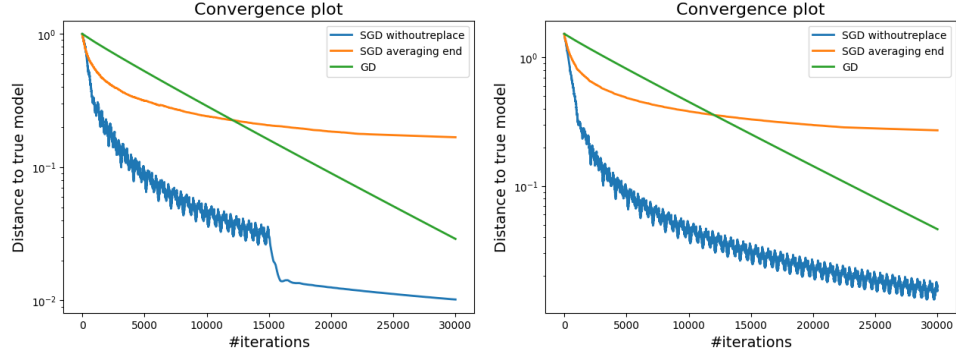


Figure 6: SGD without replacement, averaging end and GD

**How much more is the computational cost of a step of gradient descent with respect to the computational cost of a SGD step?**

A step of gradient descent (GD) has a computational cost of  $O(n)$ , while a step of stochastic gradient descent (SGD) with a single data point has a cost of  $O(1)$ . Thus, GD is  $n$  times more computationally expensive per step than SGD.

**How many steps of gradient descent should you take so that the total computational complexity is equivalent to data passes  $\times N$  steps of SGD ?**

Let  $k$  be the number steps of gradient descent:

$$k \times O(n) = \text{datapasses} \times N \times O(1) \implies k = \frac{\text{datapasses} \times N}{n} \quad (3)$$

That is the number of steps GD should take in order to get the equivalence.

**Increasing the number of data passes** helps SGD to converge to a more accurate solution. On the other hand, it leads to increased computational cost as we will need to go through the dataset multiple times.

### 3 Conclusion

In summary this implementation was amazing exploration of several techniques to improve the Stochastic Gradient Descent. Indeed those techniques really improve SGD. However, the convergence of ML algorithm does not depend only on the optimizer but also on the model architecture as well.