# Stochastic Newton with Arbitrary Sampling

**Denis Shveykin**          **Rustem Islamov**

## Abstract

The problem of minimizing the average of a large number sufficiently smooth and strongly convex functions is ubiquitous in machine learning. Stochastic first-order methods for this problem of Stochastic Gradient Descent type are well studied. In turn, second-order methods, such as Newton, have certain advances since they can adapt to the curvature of the problem. They are also known for their fast convergence rates. But stochastic variants Newton-type methods are not well-studied and have limitations on the batch size. Dmitry Kovalev et al proposed a method which requires no limitations on batch sizes. Our goal is to explore this method with different sampling strategies that lead to practical improvements.

**Keywords**

# 1 Introduction

The problem is to minimize the empirical risk which has finite-sum structure [5]:

$$\min_{x \in \mathbb{R}^d} \left[ f(x) \stackrel{def}{=} \frac{1}{n} \sum_{i=1}^{n} f_i(x) \right] \tag{1}$$

where each $f_i$ is assumed to have Lipschitz Hessian.

The stochastic approach is used, because $n$ is very large and it is computationally difficult to evaluate the gradient of each $f_i$ on each step. The Stochastic Gradient Descent (SGD) method [7] calculates the gradients of some randomly chosen $f_i$, which has cheap iterations compared to the deterministic algorithm. Such methods have two certain advantages. Firstly, they do not rely on the number of used $f_i$ on a single step, which is called the batch size. Therefore, these methods can be

applied even with small batches.

Secondly, there are ways to solve some common problem that arises due to the stochastic nature of the methods. The fact is that stochastic methods provide convergence to a neighbourhood of the solution only [6, 2]. Its size is proportional to the variance of the stochastic gradient. So the second advantage of SGD is presence of so-called variance-reduced methods [8, 4, 1, 3], which mitigate the mentioned effect. However, these methods' computational complexity depends on the curvature of the problem, which is called the condition number and is defined as the ratio of Lipschitzness and strong convexity parameters.

This leads to usage of second-order methods, such as the Newton method. Taking into account the derivatives of the second order makes it possible to adjust the algorithm's step sizes to the curvature of the problem. Unfortunately, much less literature has been written on this topic than about first-order methods. Some proposed algorithms need extra assumptions or regularization to converge. Some works provide stochastic Newton-type methods, whose computational complexity exceeds such in the variance-reduced SGD variants.

In addition, there are presented many algorithms that need large batch sizes. Particularly, the required batch size is commonly quadratically proportional to the inverse of the desired accuracy. That means that one need to evaluate a large number of $f_i$ Hessians. And this decreases the profit gained by adding randomness into the algorithm, because these batch sizes can become as big as $n$.

Dmitry Kovalev, Konstantin Mishchenko and Peter Richtarik proposed a simple Stochastic Newton algorithm, which can work with small batches, even with batches of size one. Their algorithm does not provide unbiased estimates but nevertheless shows good convergence. This is achieved by developing new Lyapunov functions that are specific to second-order methods. Our goal is to apply different sampling strategies to this method and explore their performance.

The basic set of sampling strategies can be taken from [parallel CDM], where these strategies are applied to the Parallel Coordinate Descent methods. It is also known that so-called Importance Sampling can improve the SGD performance, since it can reduce the stochastic gradient variance.

# References

[1] Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A unified theory of sgd: Variance reduction, sampling, quantization and coordinate descent, 2019.

[2] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtarik. Sgd: General analysis and improved rates. 2019.

[3] Filip Hanzely and Peter Richtárik. One method to rule them all: Variance reduction for data, parameters and many new methods, 2019.

[4] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.

[5] Dmitry Kovalev, Konstantin Mishchenko, and Peter Richtárik. Stochastic newton and cubic newton methods with simple local linear-quadratic rates, 2019.

[6] Lam Nguyen, PHUONG HA NGUYEN, Marten van Dijk, Peter Richtarik, Katya Scheinberg, and Martin Takac. SGD and hogwild! Convergence without the bounded gradients assumption. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3750–3758. PMLR, 10–15 Jul 2018.

[7] Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3):400 – 407, 1951.

[8] Nicolas Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an exponential convergence _rate for finite training sets. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.