

Cache-efficient Gradient Descent Algorithm

Imen Chakroun , Tom Vander Aa and Thomas J. Ashby
Exascience Life Lab, IMEC, Leuven , Belgium

Abstract.

Best practice when using Stochastic Gradient Descent (SGD) suggests randomising the order of training points and streaming the whole set through the learner. This results in extremely low temporal locality of access to the training set and, thus, makes minimal use of the small, fast layers of memory in an High Performance Computing (HPC) memory hierarchy. While mini-batch SGD is often used to control the noise on the gradient and make convergence smoother and more easy to identify than SGD, it suffers from the same extremely low temporal locality. In this paper we introduce Sliding Window SGD (SW-SGD) which uses temporal locality of training point access in an attempt to combine the advantages of SGD with mini batch-SGD by leveraging HPC memory hierarchies. We give initial results on a classification and a regression problems using the MNIST and ChEMBL datasets showing that memory hierarchies can be used to improve the performances of gradient algorithms.

1 Introduction

Stochastic Gradient Descent is a variety of the batch gradient descent (GD) method with a major difference in the number of updates per visited data point. In each iteration, GD sweeps through the complete training set to calculate an update to the weights vector, while for SGD only a single random element from the training data is considered. For both methods, several passes (epochs) are made over the training set until the algorithm converges or until the maximum number of iterations is reached. In a third variant called mini-batch GD (MB-GD), the model is updated based on a reasonably small groups of training samples called mini-batches. The advantages of using mini-batch gradient descent are two fold: (1) it allows the model to converge nearly as quick as SGD (in terms of time) while (2) having convergence nearly as smooth as GD. For the rest of the work MB-GD will be considered.

Since in MB-GD the search progress for the minimal cost depends on the examples picked at each iteration, best practices suggest randomizing the order of training points to visit after every epoch. This random selection is typically implemented as a random shuffling of the order of the training vectors rather than a genuine random training point selection. The shuffling is usually performed after every epoch which results in an extremely low temporal locality of access to the training set. Indeed, each training point is used once, and not before all the other training points have been visited. This means that a cache layer in the memory hierarchy of a modern HPC computer system will have little benefit for the algorithm unless all the training points fit inside that cache.

Generally speaking, a poor caching behavior of a program increasingly affects the underlying computing system's speed, cost and energy usage. Various

authors have looked at data or model parallelism for SGD to be able to benefit from the parallelism available in HPC architectures; but not on how to improve the interaction of SGD with memory hierarchies. We introduce in this paper the Sliding Window SGD (SW-SGD), a cache efficient gradient descent optimization algorithm which basic idea is to add the samples that are cached from previous iterations into the next computation of the gradient with the aim of increasing the locality of training point accesses while combining the advantages of SGD and MB-GD in terms of epoch efficiency and smoother convergence, respectively.

The remainder of this paper is as follows: in Section 2 an overview of related work is given, in Section 3 the sliding window SGD is introduced. Experimental results and associated discussion are presented in Section 4. Some conclusions and perspectives are drawn in Section 5.

2 Related work

Much research work proposing improvements on the basic stochastic gradient descent algorithm has been done. Many contributions focus on setting the learning rate because of its impact on the convergence of the algorithm. A non exhaustive list includes [2] and [1], per-parameter learning rate methods where the learning rate is adapted for each of the parameters, and per-dimension learning rates such as [3] and [4] where the step size is adapted by estimating curvature from first-order information. Some approaches to MB-GD focus on tuning the mini-batch size n such as [5] and [6, 7].

The aforementioned lines of research are more oriented towards pure algorithmic techniques compared to the contribution of this paper, which is driven by hardware considerations. In this work, the aim is to exploit hardware characteristics to improve the behaviour of mini batch SGD without incurring extra performance costs. Similar contributions considering cache locality issues in SGD but in parallel settings are [8] and [9].

3 Sliding window SGD

Today's computer architectures implement hierarchical memory structures in the attempt to alleviate the raising difference between CPU speed and main memory performance. A typical memory hierarchy contains CPU registers holding the most frequently used data, fast cache memories nearby the CPU acting as staging areas for a subset of the data and instructions stored in the relatively slow main memory and the main memory staging data stored on large, slow disks, which in turn often serve as staging areas for data stored on the disks or tapes of other machines connected by networks.

A program performance can be improved if cache blocks which have already been loaded are reused before being replaced by others. This characteristic is the motivating idea behind the SW-SGD, a cache efficient stochastic gradient descent. The gradient in SW-SGD is computed using new training points that have just been loaded from the large memory, along with some number of training

points that are still in the cache. These extra training points in the cache are essentially free to use, due to the cache effect and the fact that accessing them uses otherwise dead time whilst waiting for new points to load into the cache. Indeed, due to the higher bandwidth, the CPU can access potentially many training points in the cache in the time that it takes for a new training point to be fully loaded from the large memory into the cache. The batch size can be chosen based on performance of the resulting optimisation, and the cache size can be chosen based on hardware characteristics and performance.

Modern HPC memory hierarchies consist of many more than the two levels used to illustrate the ideas behind SW-SGD in this paper. The principle of the sliding window can be applied to every level of the memory hierarchy though, in a nested fashion, thereby extending simple SW-SGD to cover more complex hierarchies. The application of SW-SGD could be complicated somewhat by the features of the cache that is being used to store the training points. For example, hardware CPU caches have their own cache manager, replacement policy and way-mapping. A full implementation of SW-SGD for that level of memory would have to take these cache features into account to make sure that the training points were actually available in the cache as expected. By contrast, using DRAM as a cache layer for accessing a training set on disk should be relatively straight forward to implement.

From a machine learning prospective, using extra points in a gradient calculation should provide some extra smoothing similar to the extra points used in MB-GD. Ideally the smoothing would be as effective as that in MB-GD, and SW-SGD would achieve lower noise whilst having the same data touch efficiency as 1 point SGD (after accounting for pipeline-fill effects). SW-SGD differs from approaches such as momentum in that it does not introduce any more parameters to set in the algorithm. It should be possible to apply the fundamental idea of this paper to many SGD algorithmic variants (see Section 4).

4 Experiments

4.1 Experimental settings

The experiments presented in this section have been conducted on one node from a cluster of 20 each equipped with dual 6-core Intel(R) Westmere CPUs with 12 hardware threads each, a clock speed 2.80GHz and 96 GB of RAM. The implementations are all sequential (executed on one core). We used Python 3.5 and Tensorflow 1.5 as programming environment.

Two benchmark datasets have been used for the experiments. The first problem to solve is a classification problem of the MNIST dataset [10] containing 60,000 training and 10,000 testing images. The second experimented problem is a regression problem using the ChEMBL dataset which contains descriptions for biological activities involving over a million chemical entities, extracted primarily from scientific literature.

The SW-SGD have also been tested on other gradient descent optimization algorithms such as Momentum, Adam, Adagrad, etc. See [11] for a list of these

variants and an entry to the literature. The model to train is a neural network with 3 layers and 100 hidden units each. All the results are averaged from 5-fold cross-validation runs.

4.2 Experimental results

A preliminary set of experiments was conducted in order to determine the best hyper-parameters (learning rate, batch size) of the algorithm. These parameters are used for the rest of the paper.

In Figure 1, different sizes of SW-SGD are compared for different optimizers. The aim here is to first prove that SW-SGD helps accelerating the convergence and second that it is orthogonal to the other gradient algorithms. We experimented with three scenarios for every algorithm: (1) only a batch of B new points (B being the best batch size from the preliminary experiments), (2) B new points + B points from the previous iteration and (3) B new points + $2 \times B$ points from the previous iteration.

For all algorithms, adding cached data points to the computation of the gradient improves the convergence rate. For example, for the Adam gradient algorithm, a cost of 0.077 is reached after 30 epochs when using training batches of 384 points (128 new and 256 cached) while at the same epoch (30) the first scenario where no cached points are considered the cost is 0.21. Similar behavior can be observed with the other algorithms. This result proves that it should be possible to apply the fundamental idea of the SW-SGD to many GD algorithmic variants without any change to the definition of the algorithm. It is important here to highlight that, for the Adam algorithm for example, using a batch of size 256 and 512 new points is less efficient than using 128 new points as the first set of preliminary experiments showed that 128 is the best batch size. The added value here is therefore brought by the characteristic of considering old visited points in the computation and not because of a bigger batch size.

A second set of experiments was conducted on two subsets of the ChEMBL dataset. The model to train is a neural network with 3 layers and 100 hidden units each. The loss value is the root mean square error metric (RMSE). The first set we named "chembl346" has 15073 compounds and 346 proteins and the second we named "chembl360" contains 167668 compounds and 360 proteins. All the results are averaged from 5-fold cross-validation runs. The problem to solve here is to predict activity data (whether a compound is active on a protein or not) from a partially filled sparse matrix as input. The objective is to analyze the impact of the nature (sparse data) and the size of data on the performance of the SW-SGD from computing perspective in addition to the machine learning aspect. The results are presented only for SGD as gradient algorithm but the same conclusions should be drawn for the other variants.

In Table 1, the RMSE and execution times for both chembl346 and chembl360 are reported. As for the MNIST dataset, SW-SGD improves the convergence speed compared to classical mini batch SGD. From computing perspective, adding cached points to the training batch incurs extra compute time since more points are considered in the computation of the gradient. The execution

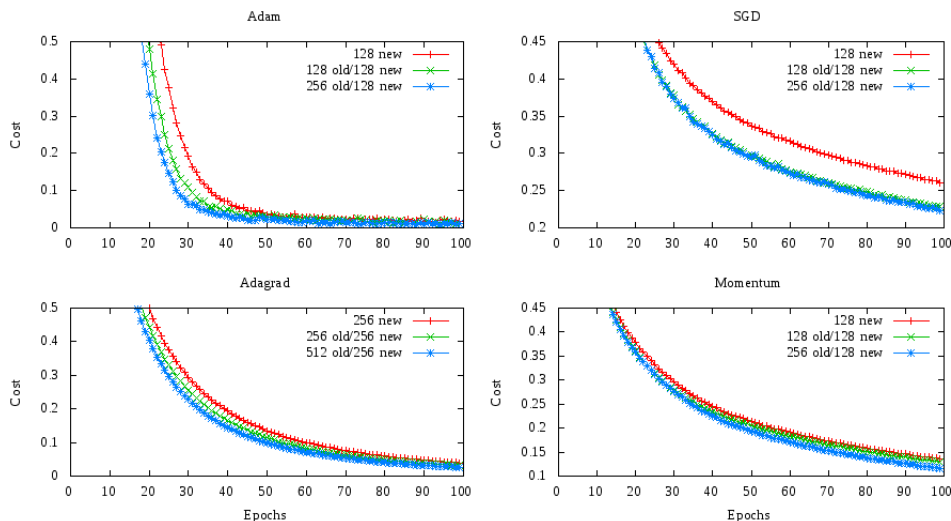


Fig. 1: Comparing different sizes of SW-SGD for different optimizers.

	chembl346			chembl360		
Batch size (new+old)	256 (256+0)	512 (256+256)	768 (256+512)	256 (256+0)	512 (256+256)	768 (256+512)
RMSE	0.364	0.27	0.228	0.67	0.563	0.509
Time in seconds	2269.69	2350.675	2418.436	52599.564	53375.83	55315.41
LL miss rate	7.7%	7.3%	7.1%	10.6%	10.1%	9.7%

Table 1: Time and RMSE for chembl346 and chembl360 using different cache sizes. The learning rate is 0.1 and the batch sizes for new points is 256.

gap is however not linear. For example, the execution time difference between batches of 256 (256 new + 0 old) and 512 (256 new + 256 old) is in the order of 3.56 % whilst the size is doubled. Moreover the bigger the size of the problem is the less the disparity is. In fact, executing 50 iterations on the chembl360 using a batch of 512 (256 new + 256 old) took 1.475 % more than when using a batch of 256 (256 new + 0 old) while it took 3.56 % more with chembl346 for the same configuration. This performance is due to the cache effect and the fact that these problems are memory and I/O bound rather than CPU bound. Most of the execution time is indeed spent in fetching and accessing the data rather than computing loss functions and gradients.

In order to highlight the effect of caching on the performances, in the third row of the Table 1, the last-level of cache (LL) miss rates are listed. These values are measured using Cachegrind which simulates how a program interacts with a machine’s cache hierarchy that exactly matches the configuration of many

modern machines. The results show that the bigger the size of the problem and the batch size are the more using SW-SGD reduces the cache hits misses.

5 Conclusion and future work

In this paper, we introduced SW-SGD which adapts mini batch GD to the access characteristics of modern HPC memory hierarchies to gain extra gradient noise smoothing, hopefully for free. In this way it will combine the epoch efficiency of one point SGD with the lower noise and easier to spot convergence of mini batch SGD. We compare the approach to mini batch SGD using different experimental data sets. We show that SW-SGD can improve over mini batch SGD in terms of convergence, for a given number of loads of training points from the large slow memory level and that is applicable for different variants of SGD.

In subsequent work we will expand the experiments to better understand under what circumstances SW-SGD should be used and how to dimension the cache. Future work should also address how to adapt the algorithm for parallel and distributed settings.

6 Acknowledgements

This work is funded by the European project ExCAPE which received funding from the European Union's Horizon 2020 Research and Innovation programme under Grant no. 671555.

References

- [1] Tieleman, Tijmen and Hinton, Geoffrey (2012). Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning
- [2] Duchi, John; Hazan, Elad; Singer, Yoram (2011). "Adaptive subgradient methods for online learning and stochastic optimization". *JMLR* 12: pp. 2121–2159.
- [3] Zeiler, Matthew D. Adadelta: An adaptive learning rate method. arXiv:1212.5701, 2012.
- [4] Schaul, T., Zhang, S., and LeCun, Y. . No more pesky learning rates. arXiv:1206.1106, 2012.
- [5] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao. Optimal distributed online prediction using mini-batches. Technical report, <http://arxiv.org/abs/1012.1367>, 2010
- [6] S. Shalev-Shwartz and T. Zhang. Accelerated mini-batch stochastic dual coordinate ascent. In *Advances in Neural Information Processing Systems*, pp. 378–385, 2013.
- [7] A. Cotter, O. Shamir, N. Srebro, and K. Sridharan. Better mini-batch algorithms via accelerated gradient methods. In *NIPS*, volume 24, pp. 1647–1655, 2011.
- [8] S. Sallinen, N. Satish, M. Smelyanskiy, S. S. Sury and C. R, "High Performance Parallel Stochastic Gradient Descent in Shared Memory," 2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS), Chicago, IL, 2016, pp. 873-882.
- [9] P. Xinghao, L. Maximilian, Tu. Stephen, Pa. Dimitris, Zh. Ce, Jo. Michael, R. Kannan, Re. Chris, Re. Benjamin, "CYCLADES: Conflict-free Asynchronous Machine Learning" . eprint arXiv:1605.09721 05/2016.
- [10] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86, 2278-2324.
- [11] <http://sebastianruder.com/optimizing-gradient-descent/index.html#gradientdescentoptimizationalgorithms>