

Block Layer Decomposition schemes for training Deep Neural Networks*

Laura Palagi, Ruggiero Seccia

Received: date / Accepted: date

Abstract Deep Feedforward Neural Networks' (DFNNs) weights estimation relies on the solution of a very large nonconvex optimization problem that may have many local (no global) minimizers, saddle points and large plateaus. As a consequence, optimization algorithms can be attracted toward local minimizers which can lead to bad solutions or can slow down the optimization process. Furthermore, the time needed to find good solutions to the training problem depends on both the number of samples and the number of variables. In this work, we show how Block Coordinate Descent (BCD) methods can be applied to improve performance of state-of-the-art algorithms by avoiding bad stationary points and flat regions. We first describe a batch BCD method able to effectively tackle the network's depth and then we further extend the algorithm proposing a *minibatch* BCD framework able to scale with respect to both the number of variables and the number of samples by embedding a BCD approach into a minibatch framework. By extensive numerical results on standard datasets for several architecture networks, we show how the application of BCD methods to the training phase of DFNNs permits to outperform standard batch and minibatch algorithms leading to an improvement on both the training phase and the generalization performance of the networks.

Keywords Deep Feedforward Neural Networks · Block coordinate decomposition · Online Optimization · Large scale optimization

* This is a preprint of an article published on Journal of Global Optimization. The final authenticated version is available online at <https://doi.org/10.1007/s10898-019-00856-0>.

The author was partially supported by the project *Distributed optimization algorithms for Big Data* of Sapienza n. RM11715C7E49E89C

Dip. di Ingegneria informatica automatica e gestionale A. Ruberti Sapienza - University of Rome
Via Ariosto 25 - 00185 Roma
Tel.: +39-06-77274081
E-mail: laura.palagi@uniroma1.it, ruggiero.seccia@uniroma1.it

1 Introduction

In this paper we consider the problem of training a Deep Feedforward Neural Network (DFNN) being available a set of training data $\{x_p, y_p\}$ for $p = 1, \dots, P$. According to the Empirical Risk Minimization (ERM) principle, training a DFNN can be formulated as an unconstrained non convex optimization problem

$$\min_{w \in \mathbb{R}^n} f(w) = \frac{1}{P} \sum_{p=1}^P E_p(w) + g(w) \quad (1)$$

where the first term represents the average loss, being E_p a measure of the loss on the single sample p , and the second term is a regularization term added to improve generalization properties which can also help from an optimization viewpoint by inducing a convexification of the objective function.

This problem is known to be a very hard optimization problem both because is highly *non-convex*, which implies the presence of stationary points that are no global minimizers, and because it is characterized by the presence of plateaus and cliffs which can extremely slow down the speed of convergence of gradient-based optimization algorithms [10,21]. This peculiar structure together with the large dimension n when considering wide and deep networks and/or the large number of samples P when considering large training set leads to a high computational effort for finding a solution of the optimization problem, being objective function and gradient evaluations the heaviest tasks in the optimization process

The optimization problem behind the training phase of DFNNs has been principally tackled with two different approaches: *batch* algorithms (a.k.a. *offline* methods), which at each iteration consider the whole dataset to update the model's weights; *minibatch* algorithms (a.k.a. *online* methods), which at each iteration consider only part of the samples in the training set, a minibatch, to update the weights of the model. The former approach can be effectively used when the number of variables n and the number of samples P in the dataset are not too large. The latter approach, by considering at each iteration only a small subset of all the available data, is more efficient when the dataset is composed by a large number of samples P . Although *minibatch* algorithms are less prone to issues when dealing with Big data, they are still affected by the dimension n of the problem, which in DFNN can grow fast.

An effective solution to solve optimization problems where the number of variables n is very large is represented by *Block Coordinate Descent* (BCD) methods [25,19,1,14], which at each iteration update only a subset of the whole variables. In contrast with *minibatch* methods, BCD methods present the opposite behaviour: they are not heavily affected by the number of variables but they still suffer when the number of samples P becomes too large.

In this paper, we study how the layered structure of DFNNs can be leveraged to define efficient BCD methods for speeding up the optimization process of these models. In particular, we introduce two different optimization frameworks:

- First, following the work done in [13] for shallow networks, we define a general *batch* BCD method for deep networks. We analyze the impact of *backpropagation* procedure on the choice of the blocks of variables and we show how a decomposition approach can help to escape from "bad" attraction regions.

- Then, exploiting both the variable decomposition, typical of BCD methods, and the sample-wise decomposition, typical of *minibatch methods* we embed the BCD scheme of above into a *minibatch* framework, defining a general *minibatch* BCD framework. Numerical results suggest how this double decomposition leads to a speed up in the solution of the non-convex optimization problem (1).

The paper is organized as follows. In Section 2 an overall literature review on both BCD algorithms, *minibatch* algorithms and how these two methods have already been combined is provided. In Section 3 some useful notation is provided and the optimization problem is formulated. In Section 4 and 5, respectively the *batch* BCD framework and the *minibatch* BCD framework are defined. In these sections, each algorithm is followed by some considerations concerning the optimization strategy implemented and numerical results on some test sets. Finally, in Section 6, conclusions and suggestions on further direction of investigations will be discussed.

2 Related works

Batch BCD algorithms are effective methods for large scale optimization problems. By updating a subset of the variables at each iteration, the application of such methods for large scale problems has already been proved to be successful in many applications [22, 19, 25, 14].

Concerning the application of BCD methods for training NNs, some approaches have already been carried out. Extreme Learning Machines (ELMs) [15, 16] are algorithms where the different role played by the output weights, namely weights to the last layer, and hidden weights are considered. Hidden weights are randomly fixed, while output weights are the tuned by an optimization procedure. The great advantage of ELMs is that when a linear unit is used in the output layer, the optimization of the output weights is a Linear Least Square problem. From the optimization point of view ELMs so not possess convergence properties to a stationary point. More sophisticated BCD methods for training Neural Networks have already been proposed. In [8] and in [13] globally convergent schemes for training respectively Radial Basis Function (RBF) NNs and MultyLayer Perceptron (MLP) were introduced where the convexity of the objective function when optimizing wrt the output weights was leveraged to improve the optimization process. These proposed frameworks have been focused on Shallow Neural Networks (SNNs), namely networks with only one hidden layer. At the best of authors' knowledge, this is the first attempt to study how a BCD scheme could be applied to deep networks which are known to be characterized by harder optimization issues.

Concerning *minibatch* algorithms, the first proposed methods for training deep networks were Incremental Gradient (IG) [3, 2] and Stochastic Gradient (SG) [23, 4, 5], where the main difference is on how minibatches are picked at each iteration, in an incremental deterministic order in IG or randomly in SG. The rate of convergence of these methods is usually slower than standard *batch* methods and depends on the variance in the gradient estimations [6]. To reduce the variance in the gradient estimation and speed up the optimization process, different approaches were implemented such as *Gradient Aggregation*, where the estimation of

the gradient is improved by considering the estimated gradients in the previous iterations. Methods like these are: SVRG [17], SAGA [11], ADAGRAD [12] and Adam [18]. As already pointed out, these methods are really efficient when the number of data is huge but are not able to scale with respect to the number of variables that for Deep Learning networks can blow up.

The behaviour of *minibatch* BCD methods has already been studied in the strongly convex case where a geometric rate of convergence in expectation has been established [24] and its effectiveness has been tested in strongly convex sparse problems such as LASSO [27] or Sparse Logistic Regression [9]. Concerning the application of *minibatch* BCD methods in training Neural Networks, in [7] a two block decomposition scheme is presented where at each iteration the output weights are exactly optimized using the full batch while the hidden weights are updated using a minibatch strategy with a step of IG.

This paper provides a general framework for the implementation of batch and online BCD methods in training DFNNs. First, filling the gap left out in [13], we study the effects of variable decomposition from both a computational point of view and an algorithmic performance point of view in *batch* algorithms for deep networks. Afterwards, we introduce a general *minibatch* BCD procedure able to speed up the training process in DFNNs and improve model's generalization properties by leveraging both the advantages of BCD and *minibatch* methods.

3 Problem Definition

Training a DFNN fits in the class of supervised learning, where a set of data $\{x_p, y_p\}_{p=1}^P$, with $x_p \in \mathbb{R}^d$ representing the input features and $y_p \in \mathbb{R}^m$ the corresponding label are given. A DFNN with d inputs and m outputs is represented by an acyclic oriented network consisting of neurons arranged in L layers connected in a feed-forward way. Each neuron j in a layer $\ell = 1, \dots, L$ is characterized by an activation function $g(\cdot)$, that for the sake of simplicity we assume the same for all the neurons, with the only exception of the output layer which has a linear activation function. We introduce the following notation that will be useful in the following.

- each layer $\ell = 1, \dots, L$ consists of N_ℓ neurons, being $N_0 = d$ and $N_L = m$ respectively the input and the output layers;
- the index set $\mathcal{L} = \{1, \dots, L\}$, where the general index ℓ is used to refer to the block of variables w_ℓ ;
- the matrix $w_\ell \in \mathbb{R}^{N_{\ell-1} \times N_\ell}$ representing the weights matrix associated to the layer $\ell \in \mathcal{L}$;
- we may refer to the variable w as composed by L blocks of variables, $w = (w_1, \dots, w_L)$;
- w_ℓ^{ji} is the weight from neuron i in layer $\ell - 1$ to neuron j in layer ℓ ;
- z_ℓ^i is the output of neuron i in layer ℓ , $z_\ell^i = g(a_\ell^i)$, where $a_\ell^i = \sum_{k=1}^{N_{\ell-1}} w_{\ell-1}^{ik} z_{\ell-1}^k$

Using this notation, given the sample x_p , we have that the output \tilde{y}_p of the DFNN can be written as:

$$\tilde{y}_p(w) = \tilde{y}(w; x_p) = w_L g(w_{L-1} g(w_{L-2} \dots g(w_1 x_p))). \quad (2)$$

We considered as error function the convex loss function MSE, and we applied a l_2 norm regularization so that the overall unconstrained problem is continuously differentiable. With these choices, the optimization problem (1) can be written as follows:

$$\min_w f(w) = \frac{1}{P} \sum_{p=1}^P E_p + \rho \|w\|^2 = \frac{1}{P} \sum_{p=1}^P \|\tilde{y}_p(w) - y_p\|^2 + \rho \|w\|^2 \quad (3)$$

As already said, this problem is highly non-convex with many local minima, saddle points and cliffs which yield to a very hard optimization problem (see [21] for a review of recent results on local- global minima issue in DNNs). Furthermore, the dependency on the number of samples and the dimension of the network make the problem even harder slowing down objective function and gradient evaluations. In this setting it is crucial to define optimization frameworks which allows to escape from bad regions and which can scale with respect to the number of samples P and the number of variables n .

4 Batch Block Coordinate Decomposition algorithm

BCD methods are iterative algorithms where, given the set of all the variables \mathcal{W} , at each iteration k only a suitable subset \mathcal{J}^k , the *working set*, is selected and only variables belonging to this set are updated by finding a new point \tilde{w}_i .

$$w_i^{k+1} = \begin{cases} w_i^k & \text{if } i \notin \mathcal{J}^k \\ \tilde{w}_i & \text{if } i \in \mathcal{J}^k \end{cases}$$

In DFNNs, at each iteration, the layered structure of the network can be exploited to select the working set. Indeed, the weights of each layer w_ℓ enters in the objective function in a nested way and a "natural" decomposition of w appears that can be fruitful used.

In this Section, taking steps from the algorithmic scheme proposed in [13] for shallow networks, we define an efficient *Batch* (BCD) scheme for the training problem of deep networks and analyze the most critical issues in its characterization. In order to show the effectiveness of the proposed choices we provide extensive numerical results on a set of benchmark datasets to point out the role of variable decomposition when dealing with deep and wide networks.

4.1 Block Layer Decomposition (BLD) scheme

We present a BCD scheme, called Block Layer Decomposition (BLD), where blocks of variables are directly defined by the layers of the network w_ℓ , $\ell = 1, \dots, L$. At the basis of this choice is the fact that fixing the weights of some layer may highlight

nice structures of the optimization problem with respect to the other variables, the working set, which makes it easier to solve (cfr [15, 13]).

In order to define the BLD algorithm, we first report in Algorithm 1 the Armijo Linesearch procedure that is used in the definition of the BLD algorithm.

Algorithm 1 Armijo Linesearch

- 1: Given $a > 0, \gamma \in (0, 1), \delta \in (0, 1)$
 - 2: Fix $\alpha = a$
 - 3: **while** $f(w^k + \alpha d^k) > f(w^k) + \gamma \alpha \nabla f(w^k)^T d^k$ **do**
 - 4: $\alpha = \delta \alpha$
 - 5: **end while**
 - 6: Return $\alpha^k = \alpha$
-

At each iteration k the proposed *batch* BLD method selects as working set an index $\ell^k \in \mathcal{L}$ which represents the block of variables belonging to layer ℓ and updates only the block w_ℓ^k by finding a point \tilde{w}_ℓ which satisfies some conditions. Then a new index ℓ_{k+1} is extracted according to some rule and the same steps are repeated until a stopping criterion, such as a maximum number of iterations or the norm of the gradient below a certain threshold, is met.

In particular, the condition which must be satisfied by the new point block \tilde{w}_ℓ are:

1. The objective function evaluated at the new point $w^{k+1} = (w_1, \dots, \tilde{w}_\ell, \dots, w_L)$ is not worse than the value obtained along the steepest descent direction $d_\ell^k = -\nabla_{w_\ell} f(w^k)$ with the stepsize chosen by an Armijo Linesearch

$$f(w_1^k, \dots, \tilde{w}_\ell, \dots, w_\ell^k) \leq f(w_1^k, \dots, w_\ell^k + \alpha^k d_\ell^k, \dots, w_\ell^k) \quad (4)$$

2. The difference between the objective function evaluated in the point w^k and in the new one satisfies

$$f(w_1^k, \dots, \tilde{w}_\ell, \dots, w_\ell^k) - f(w_1^k, \dots, w_\ell^k, \dots, w_\ell^k) \leq -\sigma(\|\tilde{w}_\ell - w_\ell^k\|) \quad (5)$$

where σ is a forcing function satisfying

$$\lim_{k \rightarrow \infty} \sigma(t_k) = 0 \implies \lim_{k \rightarrow \infty} t_k = 0$$

The Batch BLD scheme is reported in Algorithm 2.

Concerning convergence to stationary non-maxima points of the Batch BLD framework, we have that it fits within the general decomposition scheme proposed in Section 7 of [14]. To prove convergence we need to introduce the following assumption on the selection of the index set ℓ^k :

Assumption 1 (Cyclic updating rule) *Blocks of variables w_ℓ must be updated in a cyclical order*

Under this assumption, we can state the following convergence result whose proof, for the sake of completeness, is reported in the Appendix.

Theorem 1 *Algorithm BLD with a cyclical selection of the blocks generates a sequence of points $\{w^k\}$ such that*

$$\lim_{k \rightarrow \infty} \nabla f(w^k) = 0 \quad (6)$$

Algorithm 2 *Batch* Block Layer Decomposition scheme

```

1: Given  $\{x_p, y_p\}_{p=1}^P, \mathcal{L} = \{1, \dots, L\}$ 
2: Choose  $w^0 \in \mathbb{R}^n$  and set  $k = 0$ ;
3: while (stopping criterion not met) do
4:   Select an index  $\ell^k \subseteq \mathcal{L}$ 
5:   Set  $d_{\ell^k} = -\nabla_{w_{\ell^k}} f(w^k)$ 
6:   if  $d_{\ell^k} \neq 0$  then
7:     Compute  $\alpha^k$  along  $d_{\ell^k}$  with an Armijo Linesearch
8:     Find a point  $\tilde{w}_{\ell^k}$  such that (4) and (5) are satisfied
9:   else
10:     $\tilde{w}_{\ell^k} = w_{\ell^k}^k$ 
11:   end if
12:   Update  $w^{k+1} = (w_1^k, \dots, \tilde{w}_{\ell^k}, \dots, w_L^k)$ 
13:   Set  $k = k + 1$ 
14: end while

```

Algorithm 2 must be further characterized depending on how the index ℓ^k is chosen and how the new point \tilde{w}_ℓ is computed. These choices affect both the convergence properties of the algorithm and its computational efficiency. In the following, we discuss the possible choices that satisfy assumptions to guarantee convergence of BLD and allow to save computations, mainly due to the reduction of the needed time for objective function and gradient evaluations.

Selection of the *working set* ℓ^k

Before providing general rules on how to choose the index ℓ^k , we need to recall how the *backpropagation* procedure works in order to comprehend the crucial role played by the working set definition in the BLD scheme.

In DFNNs, the main effort in computing $\nabla f(w)$ stays in evaluating $\sum_{p=1}^P \nabla_w E_p$. Thanks to the chain rule, each element $\nabla_w E_p$ is computed as

$$\frac{\partial E_p(w)}{\partial w_\ell^{ji}} = \frac{2}{P} z_{\ell-1}^i \delta_\ell^j \quad (7)$$

where the z_ℓ are obtained by forward propagation of the input and the δ_ℓ by backward propagation throughout the layers of the error $e = \tilde{y}_p - y_p$

$$\delta_L^j = e g'(a_L) \quad \delta_\ell^j = \sum_{k=1}^{N_{\ell+1}} \delta_{\ell+1}^k w_{\ell+1}^{kj} g'(a_\ell^j) \quad l \leq L - 1 \quad (8)$$

where g' is the partial derivative of g .

This procedure is called *backpropagation* because, to get derivatives of variables in layer ℓ , we need to evaluate the δ s of all the following layers, starting from the last one δ_L and backpropagating the δ s until δ_ℓ , being the general δ_ℓ depending on $\delta_{\ell+1}$. Then, if the full gradient $\nabla f(w^k)$ must be evaluated, as in standard gradient based methods, forward and backward propagation involves all the layers of the network. On the other hand, when only the gradient wrt a block ℓ is needed, as in BLD scheme, the propagation steps go only over *part* of the network, reducing the computational effort.

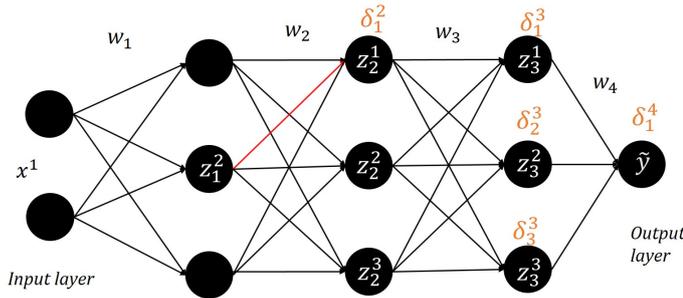


Fig. 1 Example of how *backpropagation* works to get the derivative wrt a hidden weight

In order to save computations in the evaluation of the gradient, the chosen block layer decomposition turns out to be the more efficient. Indeed, looking at (7) and (8), it is clear that blocks of variables defined with interlacing weights over the layers would imply the forward and backward propagation along the full network for the computations of δ s. Furthermore, we highlight that also a definition of blocks of variables by neurons, as done in [13] for shallow networks, turns out to be inefficient, since to update the general neuron j at layer ℓ we need to compute δ_ℓ which depends by all the δ s of the following layers which would be computed but not further used. Another possible choice for exploiting the backpropagation procedure consists in updating at iteration k the weights of layer ℓ^k and all the successive layers $\ell^k + 1, \dots, L$, being the following δ s already computed to evaluate δ_ℓ . This blocks selection rule would be more efficient since it would imply the utilization of all the δ s computed at each iteration. However, this procedure leads to an unbalanced decomposition update because the last layers are optimized more often than the first ones. Extensive numerical results (not reported here for the lack of space) shows a worst performance as final value of the objective function. This behaviour is arguably due to the fact that this strategy gives more importance to some layers than others making the algorithm converge to very poor regions of the problem and severely harming the overall training process.

As regard the potential strategies for choosing the index ℓ^k , Assumption (1) can be ensured by using several block selection rules such as:

- **Forward selection.** Layers are selected according to the forward propagation $\{1, 2, \dots, L\}$.
- **Backward selection.** layers are selected according to the backward propagation $\{L, L - 1, \dots, 1\}$.
- **Random without replacement.** The blocks are selected with a cyclic order that changes after all the blocks have been updated. This corresponds to a cyclic rule with random reshuffling. In this case

Block \tilde{w}_ℓ update

Once block ℓ^k is selected, a point \tilde{w}_{ℓ^k} satisfying conditions (4) (5) must be determined. First we note that, similarly to Extreme Learning Machine [15,16], when the last block $\ell = L$ is selected, the optimization problem becomes a linear least square problem (LLSQ) with an additional $l-2$ regularization term which makes

the subproblem in the variables w_L strictly convex. As a consequence, an update \tilde{w}_L satisfying conditions 4, 5 can be easily found either by solving the subproblem in closed form or by the application of an iterative method such as the Conjugate Gradient algorithm [20]. Note that since the optimization wrt w_L is quadratic strictly convex, the update $\tilde{w}_L^{k+1} = \arg \min_{w_L} f(w^k)$ is always bounded to guarantee conditions (4)(5) (see [14]).

When a block different from the last one is selected, $\ell \neq L$, the optimization problem becomes non-convex and finding good updates \tilde{w}_ℓ may become harder. We first note that the point returned by the Armijo Linesearch along the steepest descent direction

$$\tilde{w}_{\ell^k} = w_{\ell^k}^k - \alpha^k \nabla_{w_{\ell^k}} f(w^k)$$

satisfies the two conditions (4) (5). However a single iteration of the gradient method for each ℓ^k can slow down the overall procedure. More in general, it is easy to prove that a finite number of steps of a gradient method with Armijo Linesearch guarantees conditions (4) (5) (see Appendix for more details), so that \tilde{w}_{ℓ^k} can be obtained by iterating a gradient method with an Armijo Linesearch. However, as suggested in [13], since the number of variables in NNs may be very high, methods which approximate second-order information and do not suffer too much of the *curse of dimensionality* would perform remarkably better than gradient methods. So, at each iteration an optimization method such as a limited memory Quasi-Newton method [20] like LBFGS can be applied to solve the subproblem in the variables w_{ℓ^k} to get a trial point $w_{\ell^k}^*$. If conditions 4, 5 are satisfied, $\tilde{w}_{\ell^k} = w_{\ell^k}^*$, otherwise \tilde{w}_{ℓ^k} is set to the point obtained by an Armijo Linesearch along the steepest descent direction.

4.2 Implementation and performance of a BLD method

For the experimental results, we implemented a BLD algorithm with the following settings:

- a *backward cyclic order* selection rule of the layer;
- choice of \tilde{w}_ℓ^k by the application of a LBFGS method with increasing accuracy ε and limited number of iterations with a check on conditions (4) (5) as described above.

We observe that we applied an LBFGS method also to solve the LLSQ problem in the variables of the last block w_L . This choice is due to the observation in the numerical results that forcing a high accuracy in the optimization of the last layer from the early stages seemed to get worse results in terms of quality of the solution (value of the objective function and test error as well).

We compared B²LD performance against an LBFGS method applied to the whole optimization problem 3. The initial point w^0 has been set to the same randomly chosen value for both the two methods and performance are compared on the results obtained within 10 runs. The regularization parameter ρ was set equal to $\rho = 10^{-3}/n$ and the sigmoid function was used as activation function. No hyperparameters tuning was carried since the focus is on studying optimization algorithms performance more than finding the smaller generalization error.

Algorithm 3 Backward Block Layer Decomposition (B²LD)

```

1: Given  $\{x_p, y_p\}_{p=1}^P$ 
2: Choose  $w^0 \in \mathbb{R}^n$  and set  $k = 0$ ; set  $\epsilon > 0, \delta \in (0, 1)$ .
3: while stopping criterion not met do
4:   for  $l = L, \dots, 1$  do
5:     Find  $\tilde{w}_\ell$  s.t.  $\|\nabla_{w_\ell} f(\tilde{w}_\ell)\| \leq \epsilon$ 
6:     if  $\tilde{w}_\ell$  satisfies (4) (5) then
7:        $w^{k+1} = (w_1^k, \dots, \tilde{w}_\ell, \dots, w_L^k)$ 
8:     else
9:        $w^{k+1} = (w_1^k, \dots, w_\ell^k - \alpha^k \nabla_{w_\ell} f(w^k), \dots, w_L^k)$ 
10:    end if
11:     $k = k + 1$ 
12:  end for
13:   $\epsilon = \epsilon \delta$ 
14: end while

```

The following stopping criteria were set: i) $\|\nabla f(w^k)\| \leq 10^{-3}$; ii) $f_{tol} = \frac{f(w^k) - f(w^{k+1})}{\max\{f(w^k), 1\}} \leq 10^{-4}$; iii) computational time exceeds 150 seconds. Note that for B²LD condition ii) must be satisfied with respect to each block in order to stop computations. Moreover, computational times were checked only every 30 iterations of LBFGS method that is why in Table 2 Cpu Times can be higher than 150 seconds.

Both the two algorithms were implemented in Python version 3.6, using the package *numpy* for numerical calculus and *scipy* for the optimization algorithms using a Intel(R) Core(TM) i7-870 CPU 2.93 GHz.

The two algorithms have been compared over five different network architectures in order to analyze how the algorithms perform when dimensions of the problems increase with respect to two different parameters: width and depth of the network. The characteristic of the different architectures are reported in Table 1. For the sake of simplicity in most cases we assume equal numbers of neurons per layer $N_\ell = N$, so that each network is identified by the notation $[L \times N]$, with the exception of the network $[200, 50, 200]$, where the values represent the number of neurons per the three layers. Comparison were carried over seven different datasets publicly available at <https://sci2s.ugr.es/keel/index.php> and <https://archive.ics.uci.edu/ml/index.php>. The number of samples in the training and test sets and the number of features per dataset are reported in Table 1 along with the number of variables per each optimization problem.

Table 1 Datasets and Networks Description

Dataset	# Train	# Test	# Features	# Variables in the network				
				[1 x 50]	[3 x 20]	[200,50,200]	[5 x 50]	[10 x 50]
Mv	32614	8154	13	700	1080	22800	10700	23200
Bikes Sharing	13903	3476	59	3000	2000	32000	13000	25500
Beijing PM 2.5	33406	8351	48	2450	1780	29800	12450	24950
CCPP	7654	1913	5	300	920	21200	10300	22800
Ailerons	11000	2750	41	2100	1640	28400	12100	24600
California	16512	4128	9	500	1000	22000	10500	23000
Bank	36168	9042	40	2050	1620	28200	12050	24550

These numerical experiments aim to show that using BLD may help in obtaining better results in term of the final objective function (regularized training error) without deteriorating the generalization performance. We analyze both the best performance over the 10 runs and the average behaviour.

In Table 2 the detailed results corresponding to the best run of each of the two algorithms are reported. For each network architecture and for each problem we report: the best solution found, the corresponding norm of the gradient and the time needed to find it.

Analyzing the numerical results, we observe that BLD method is able to find better values of the objective function with a smaller value of the norm of the gradient. It seems that BLD does not get stuck in poor solutions as it happens to the standard LBFGS method. This issue is particularly evident when the number of layers increases. Indeed, the results on the deepest network $[10 \times 50]$ show that LBFGS stops very quickly at a very poor solution with a small norm of the gradient whereas BLD is able to continue optimization. This may be due also the *vanishing gradient* effect which may affect a full gradient method.

In order to get a quick look at the results we also produce histograms in which we analyze the performance over each problem by changing the architectures. In this way it is possible to visualize how the performance over the same problem changes with the architecture of the network. B²LD's performance seems to be not affected by the number of layers, whereas for deeper network LBFGS seems to be trapped into "bad" solutions.

The average behaviour of B²LD is assessed by counting the number of win, #wins, (B²LD beats LBFGS) and defeats, #defs, (LBFGS beats B²LD) over the 10 random runs. We say that a method beats the other if the returned value is better than 5% compared with the other; otherwise we declare a tie. The cumulative results are reported in Table 3 and 4 respectively on the objective function and the test error value. For each problem and for each network architecture we report into parenthesis [#wins ; #defs] of B²LD w.r.t. LBFG. From this analysis it is evident that B²LD remarkably outperforms LBFGS particularly when the dimension of the problem increases; indeed it obtains better solutions (smaller objective function, Table 3) with better generalization properties (smaller test error, Table 4).

Finally, we analyzed how many times each layer is updated when using B²LD. Indeed, while LBFGS method updates the weights of all the layers at the same time, B²LD considers layers one by one and it is possible that some layer is updated less often than others. Indeed, whenever the gradient with respect to a layer and/or the relative reduction of the objective function are below a given tolerance, B²LD skips the update of the layer. In this way, the algorithm does not waste time in optimizing parts of the network not worthwhile. We consider the deepest network, $[10 \times 50]$, and we count the number of updates performed by the best BLD on each layer, which are reported in Table 5. It turns out that the first and last layers are optimized much often than layers in the middle of the network. This highlight the fact that the network might be too deep for the problems and a smaller number of hidden layers may be enough. We also note that on the dataset Beijing PM 2.5, B²LD performs only the optimization of the last layer, miming the behaviour of Extreme Learning Machine and it finds a better solution than LBFGS (see Table 2).

Table 2 Best result returned over 10 runs by the B²LD and LBFGS methods

Network	Objective Function		Norm of the Gradient		Cpu Time	
	B ² LD	LBFGS	B ² LD	LBFGS	B ² LD	LBFGS
[1 x 50]						
Ailerons	2.37×10^{-3}	2.32×10^{-3}	9.98×10^{-4}	4.03×10^{-4}	0.18	3.04
Bank Marketing	7.92×10^{-5}	3.67×10^{-5}	4.96×10^{-4}	6.06×10^{-4}	13.91	29.77
Bejing Pm25	4.79×10^{-3}	5.32×10^{-3}	8.57×10^{-4}	2.25×10^{-3}	9.16	8.43
Bikes Sharing	2.04×10^{-3}	3.01×10^{-3}	7.12×10^{-4}	4.10×10^{-3}	16.44	11.58
California	2.00×10^{-2}	2.00×10^{-2}	2.98×10^{-4}	3.34×10^{-3}	1.50	2.99
CCPP	3.35×10^{-3}	3.53×10^{-3}	5.88×10^{-4}	6.87×10^{-4}	0.08	0.91
Mv	2.93×10^{-4}	2.68×10^{-4}	7.18×10^{-4}	4.90×10^{-4}	20.90	23.78
[3 x 20]						
Ailerons	2.33×10^{-3}	2.24×10^{-3}	9.29×10^{-4}	4.88×10^{-4}	3.50	3.62
Bank Marketing	1.09×10^{-4}	3.91×10^{-5}	8.09×10^{-4}	7.66×10^{-4}	59.19	17.45
Bejing Pm25	4.45×10^{-3}	7.60×10^{-3}	9.48×10^{-4}	1.36×10^{-3}	6.24	6.04
Bikes Sharing	1.88×10^{-3}	2.81×10^{-3}	8.37×10^{-4}	2.36×10^{-3}	10.76	16.54
California	2.02×10^{-2}	2.25×10^{-2}	4.10×10^{-4}	2.28×10^{-3}	2.44	1.56
CCPP	3.47×10^{-3}	3.35×10^{-3}	3.89×10^{-4}	3.38×10^{-4}	1.51	1.17
Mv	2.92×10^{-4}	2.17×10^{-4}	7.71×10^{-4}	3.70×10^{-4}	23.86	37.44
[200.50.200]						
Ailerons	2.33×10^{-3}	2.56×10^{-3}	9.09×10^{-4}	1.58×10^{-3}	17.93	36.77
Bank Marketing	1.93×10^{-4}	2.31×10^{-4}	1.61×10^{-3}	4.98×10^{-3}	184.34	151.09
Bejing Pm25	4.59×10^{-3}	5.76×10^{-3}	8.44×10^{-4}	2.59×10^{-3}	137.09	113.77
Bikes Sharing	2.01×10^{-3}	4.00×10^{-3}	5.81×10^{-4}	3.12×10^{-3}	146.89	132.06
California	1.98×10^{-2}	2.31×10^{-2}	1.14×10^{-2}	5.46×10^{-3}	17.37	15.41
CCPP	3.50×10^{-3}	3.70×10^{-3}	7.04×10^{-4}	8.17×10^{-4}	6.67	16.39
Mv	2.49×10^{-4}	9.73×10^{-4}	7.84×10^{-4}	5.59×10^{-3}	145.32	185.38
[5 x 50]						
Ailerons	2.37×10^{-3}	4.89×10^{-3}	8.15×10^{-4}	5.49×10^{-3}	12.57	10.52
Bank Marketing	2.02×10^{-4}	9.19×10^{-5}	9.41×10^{-4}	1.43×10^{-3}	110.51	152.88
Bejing Pm25	5.03×10^{-3}	7.26×10^{-3}	5.38×10^{-4}	1.34×10^{-3}	42.72	27.26
Bikes Sharing	2.07×10^{-3}	3.46×10^{-2}	5.70×10^{-4}	5.06×10^{-3}	74.30	1.81
California	2.05×10^{-2}	5.47×10^{-2}	2.92×10^{-3}	1.57×10^{-3}	12.56	1.46
CCPP	3.79×10^{-3}	4.45×10^{-3}	4.54×10^{-4}	1.54×10^{-3}	7.33	8.61
Mv	3.59×10^{-4}	7.76×10^{-4}	6.12×10^{-4}	8.38×10^{-3}	131.44	150.55
[10 x 50]						
Ailerons	2.65×10^{-3}	1.31×10^{-2}	2.92×10^{-2}	4.11×10^{-5}	68.38	1.95
Bank Marketing	2.10×10^{-3}	6.28×10^{-2}	1.94×10^{-01}	7.93×10^{-5}	188.25	6.77
Bejing Pm25	8.59×10^{-3}	8.74×10^{-3}	8.18×10^{-4}	1.77×10^{-5}	0.21	5.08
Bikes Sharing	2.68×10^{-3}	3.51×10^{-2}	3.34×10^{-2}	6.89×10^{-5}	158.46	2.65
California	1.88×10^{-2}	5.50×10^{-2}	1.31×10^{-01}	1.35×10^{-4}	75.67	2.65
CCPP	4.89×10^{-3}	5.16×10^{-2}	1.61×10^{-2}	1.12×10^{-4}	17.05	1.16
Mv	3.20×10^{-3}	5.53×10^{-2}	1.76×10^{-01}	1.16×10^{-4}	167.36	4.85

Table 3 [#wins ; #defs] of B²LD versus LBFGS on the objective function values (10 runs)

	[1 x 50]	[3 x 20]	[200, 50, 200]	[5 x 50]	[10 x 50]
Ailerons	[0 ; 0]	[6 ; 2]	[8 ; 0]	[10 ; 0]	[9 ; 0]
Bank Marketing	[0 ; 10]	[5 ; 5]	[4 ; 5]	[8 ; 2]	[10 ; 0]
Bejing Pm25	[5 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]	[0 ; 0]
Bikes Sharing	[10 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]
California	[8 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]
CCPP	[2 ; 0]	[2 ; 6]	[8 ; 0]	[10 ; 0]	[10 ; 0]
Mv	[3 ; 5]	[5 ; 3]	[10 ; 0]	[10 ; 0]	[10 ; 0]

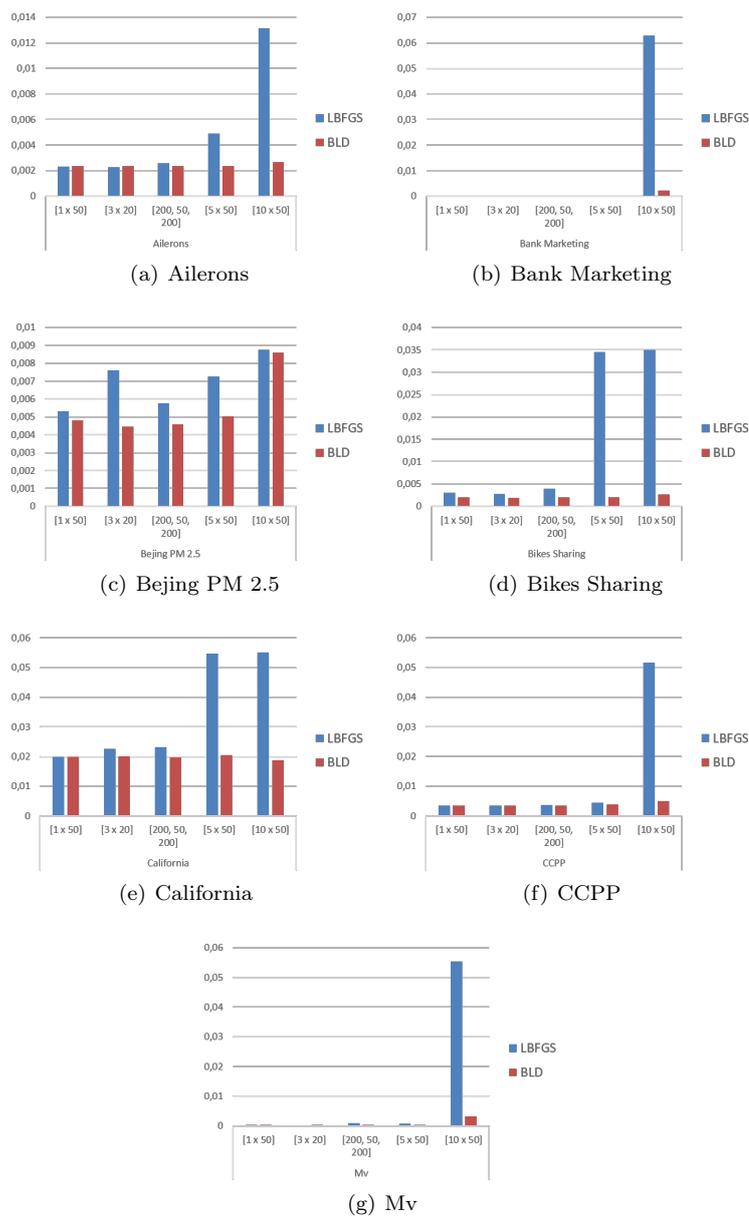


Fig. 2 Best Objective function value returned by the BLD and the LBFSGS methods.

Table 4 [#wins ; #defs] of B²LD versus LBFSGS on the test error values (10 runs)

	[1 x 50]	[3 x 20]	[200, 50, 200]	[5 x 50]	[10 x 50]
Ailerons	[0 ; 0]	[6 ; 4]	[7 ; 0]	[10 ; 0]	[9 ; 0]
Bank Marketing	[0 ; 9]	[5 ; 5]	[2 ; 8]	[8 ; 2]	[10 ; 0]
Bejing Pm25	[5 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]	[0 ; 0]
Bikes Sharing	[10 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]
California	[8 ; 0]	[9 ; 0]	[10 ; 0]	[10 ; 0]	[10 ; 0]
CCPP	[0 ; 0]	[1 ; 8]	[8 ; 0]	[9 ; 1]	[10 ; 0]
Mv	[3 ; 6]	[4 ; 6]	[10 ; 0]	[8 ; 1]	[10 ; 0]

Table 5 Number of updates per layer $\ell = 0, \dots, 10$ performed by B²LD during the best run for the network [10 x 50].

Dataset	Layer of the network										
	0	1	2	3	4	5	6	7	8	9	10
Mv	30	30	48	8	2	2	2	2	2	2	14
Bikes Sharing	120	62	35	6	62	5	4	4	4	4	18
Bejing PM 2.5	0	0	0	0	0	0	0	0	0	0	5
CCPP	11	2	2	31	2	2	2	2	2	2	4
Ailerons	33	16	24	63	4	7	4	4	4	4	21
California	28	28	39	11	2	2	2	2	2	2	13
Bank Marketing	30	30	30	30	1	6	2	30	1	1	5

5 Minibatch Block Layer Decomposition algorithm

So far, we have seen that the BLD method can improve the performance of standard optimization algorithms. In this section, we exploited the additive structure of the objective function and we embed the layer decomposition scheme proposed above into a *minibatch* strategy with the aim of enhancing the performance of both online methods and batch decomposition methods. A similar idea has been proposed in [7] where a two-block incremental decomposition method has been proposed which is suitable for block layer decomposition of a shallow network.

5.1 Minibatch Block Decomposition Method

Let $\mathcal{B}_h \subset \{1, \dots, P\}$ be the index set identifying a minibatch, i.e. a subset of the samples. Given a partition $\mathcal{B} = \{\mathcal{B}_1, \dots, \mathcal{B}_H\}$ of the set $\{1, \dots, P\}$, problem (3) can be expressed as

$$\min_w \sum_{h=1}^H \sum_{p \in \mathcal{B}_h} E_p(w_1, \dots, w_L) + \rho \|w\|^2$$

The *minibatch* BCD algorithm applies a double decomposition to the problem, meaning that at each iteration only information on the function involving a minibatch \mathcal{B}_h is used to update a subset of variables, the *working set*, \mathcal{J}^k .

Let us introduce the following notation

$$f_h = \sum_{p \in \mathcal{B}_h} E_p(w_1, \dots, w_L) + |\mathcal{B}_h| \frac{\rho}{P} \|w\|^2$$

so that $f(w) = \sum_{h=1}^H f_h$ and the gradient can be expressed as

$$\nabla f(w) = \sum_{h=1}^H \sum_{\ell=1}^L e_{\ell} \otimes \nabla_{w_{\ell}} f_h$$

where \otimes represents the standard Kronecker Product and $e_{\ell} \in \mathbf{R}^L$ is the ℓ -th column of the identity matrix of dimension L .

A very general framework of a minibatch block layer decomposition is reported in Algorithm 4.

Algorithm 4 Minibatch Block Layer Decomposition

```

1: Given  $\{x_p, y_p\}_{p=1}^P$ 
2: Choose  $w^0 \in \mathbb{R}^n$  and set  $k = 0$ ;
3: while (stopping criterion not met) do
4:   Define a partition  $\mathcal{B}^k = \{\mathcal{B}_1, \dots, \mathcal{B}_H\}$  of  $\{1, \dots, P\}$ 
5:   Select a minibatch  $\mathcal{B}_h, h \in \{1, \dots, H\}$ 
6:   Choose  $\mathcal{J}_h^k \subseteq \{1, \dots, L\}$ 
7:   Set  $\tilde{w}^0 = w^k; j = 0$ ;
8:   for  $\ell \in \mathcal{J}_h^k$  do
9:     Set  $d_{\ell}^j = -\nabla_{w_{\ell}} f_h(\tilde{w}^j)$ 
10:    Update  $\tilde{w}_{\ell}^{j+1} = \tilde{w}_{\ell}^j + \alpha^k d_{\ell}^j$ 
11:     $j = j + 1$ 
12:  end for
13:   $w^{k+1} = \tilde{w}^{|\mathcal{J}_h^k|}$ 
14:  Update  $\alpha^k$ 
15:   $k = k + 1$ 
16: end while

```

Algorithm 4 can be differently characterized depending on the definition of the partition \mathcal{B}^k , the selection rule of the minibatch \mathcal{B}_h , the choice of the working set \mathcal{J}_h^k , and the updating rule of the stepsize α^k . We address possible interesting choices in the following paragraphs.

Selection of the minibatch \mathcal{B}_h

Selection of the minibatch over the P samples can be implemented following classical rules used in *online* methods for Machine Learning. In particular we can consider

- **Incremental Rule:** the order in which minibatches will be used is fixed *a priori* and kept unchanged over the iterations;
- **Stochastic Rule:** at each iteration a minibatch is chosen randomly from the available list \mathcal{B} ;
- **Random without replacement rule:** at each iteration a minibatch is chosen randomly in \mathcal{B} without replacement; this corresponds to an incremental rule when the selection order is reshuffled;

Selection of the *working set* \mathcal{J}_h^k

The index set \mathcal{J}_h^k defines which blocks will be sequentially updated at iteration k using the minibatch \mathcal{B}_h . \mathcal{J}_h^k can be composed by all the layers, $\mathcal{J}_h^k = \{1, \dots, L\}$, or can be a subset of the index set \mathcal{L} . In the first case, all the blocks are updated sequentially using the same minibatch, while, in the second case, only a subset of the layers are updated with a given minibatch. Furthermore, \mathcal{J}_h^k can be defined at each iteration, or it can be fixed *a priori*, $\mathcal{J}_h^k = \mathcal{J} \forall k$.

In the definition of the set \mathcal{J}^k , a computationally efficient choice, that exploits the backpropagation rule for updating the gradient, consists in setting set $\mathcal{J}^k = \{1, \dots, L\}$. Indeed, when updating a block of variables w_{ℓ^k} with a minibatch \mathcal{B}_h in each layer ℓ we can store the output z_ℓ^i of each neuron i so that if the same minibatch is used to update another block of variables $w_{\ell^{k+1}}$ we do not need to evaluate the output z_ℓ of the previous layers $\ell < \ell^{k+1}$, being these unchanged.

5.2 Implementation and performance of a minibatch BLD method

For the experimental results, we implemented a *minibatch* BCD algorithm with the following settings:

- the order in which minibatches \mathcal{B}_h are used is fixed *a priori*;
- at each iteration the set \mathcal{J}_h^k is composed by all the layers which are updated sequentially following a *backward* order as in the B²LD scheme presented above;
- α^k updated according to the diminishing stepsize

$$\alpha^k(1 - \epsilon\alpha^k)$$

At each iteration, the algorithm picks a minibatch \mathcal{B}_h and updates the layers of the network sequentially in backward order by performing a steepest descent iteration. After having updated all the layers, a new minibatch is considered and the same procedure is applied until a certain stopping criterion, such as a maximum number of iterations or the norm of the gradient smaller than a certain threshold, is met. We call this scheme Block Layer Incremental Gradient (BLInG) which is reported in Algorithm 5.

Convergence of the BLInG above can be proved under suitable assumptions following [4] by looking at the iteration generated by BLInG as a gradient method with error. A similar approach has been followed in [7] for a two-block decomposition where the last block is optimized with a full batch strategy.

Following [26], in the updating formula of w_ℓ^k we have scaled the stepsize with a normalization term, properly bounded to avoid overflow, since this choice seemed to make the updating process more robust by avoiding vanishing and exploding gradient.

As regards the updating rule for the stepsize α^k , convergence requires to use a *diminishing* stepsize rule. The proposed rule at Step 13 of the BLInG algorithm satisfies the assumptions as well as another usual updating rule $\alpha^{k+1} = \frac{1}{k}\alpha^k$. Concerning the parameters in the updating rule of the stepsize, γ and β were fixed to 10^{-3} and 10^6 respectively.

We compared BLInG with the Incremental Gradient (IG), which is its non-decomposed counterpart. We use for both the two algorithms the same setting and

Algorithm 5 Block Layer Incremental Gradient (BLInG)

```

1: Given  $\{x_p, y_p\}_{p=1}^P, \mathcal{L} = \{1, \dots, L\}$ 
2: Choose  $w^0 \in \mathbb{R}^n$  and set  $k = 0, \alpha^0 \geq 0, \epsilon \in (0; 1), \beta > 0, \gamma < \infty;$ 
3: Define a partition  $\mathcal{B} = \{\mathcal{B}_1, \dots, \mathcal{B}_H\}$  of  $\{1, \dots, P\}$ 
4: while (stopping criterion not met) do
5:   for  $h=1, \dots, H$  do
6:     Set  $\tilde{w}^0 = w^k; j = 0;$ 
7:     for  $\ell = 1, \dots, L$  do
8:       Set  $d_\ell^j = -\nabla_{w_\ell} f_h(\tilde{w}^j)$ 
9:       Update  $\tilde{w}_\ell^{j+1} = \tilde{w}_\ell^j + \frac{\alpha^k}{\max\{\beta, \min\{\gamma, \|d_\ell^k\|\}\}} d_\ell^j$ 
10:       $j = j + 1$ 
11:    end for
12:     $w^{k+1} = \tilde{w}^{|\mathcal{J}_h^k|}$ 
13:     $\alpha^{k+1} = \alpha^k (1 - \epsilon \alpha^k)$ 
14:     $k = k + 1$ 
15:  end for
16: end while

```

normalization for the stepsize α^k . The initial values of the learning rate α^0 and the fraction reduction ϵ were chosen through a grid-search procedure which led to different values for the two methods which depend on the network architecture:

$$\alpha_{IG}^0 = 0.5 \quad \alpha_{BLInG}^0 = \frac{0.5}{\max\{1, L - 2\}} \quad \epsilon = 5 \times 10^{-3}$$

IG performance was invariant with respect to the initial value of the learning rate which best value was the same for all the architectures of the network; BLInG instead performs better with a smaller initial learning rate for deeper networks.

Datasets and architectures of the networks are the same of those used for the BLD algorithm and reported in Table 1. Each algorithm was tested over 10 runs starting from randomly chosen initial points.

Since *minibatch* methods evaluate neither the objective function nor the gradient at each iteration, the stopping criterion has been a limit on the computational time that was fixed to 60 seconds.

The results are reported in Tables 6 and 7. In Table 6 for each of the two algorithms, we report the best values of the objective function found over the 10 runs and the corresponding value of the test error within the time limit of 60 sec. Overall, BLInG is able to return better solutions with better generalization properties on the test set than those returned by IG, especially when the dimension of the network blows up. Table 7, provides cumulative information on the 10 runs. In particular, we considered the number of wins or defeats of the BLInG algorithm versus IG over the 10 runs (being a tie a result within the 5%). Similarly to the BLD algorithm, BLInG seems to perform better when the dimension of the problem increases. In order to assess how much the depth of the network influences algorithms' performance, in Table 8 the ratio between best value found in the network [10x50] and [1x50] are provided. Higher values mean that the algorithm has been harmed by the increased depth of the network. Comparing these solutions, BLInG turns out to be less affected by the increased structure of the network and is able to find always similar values regardless of the structure of the network while IG performs worse.

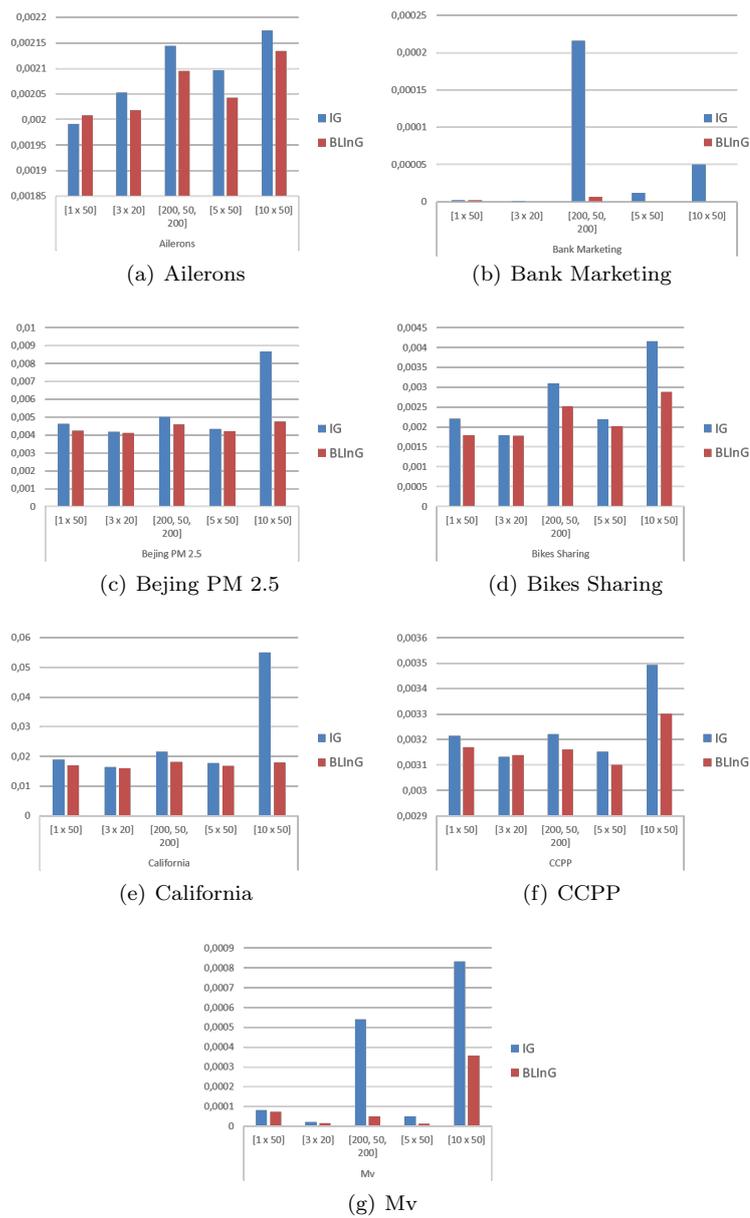


Fig. 3 Best training error value obtained by IG and BLInG within 60 seconds.

Table 6 Best results returned over 10 runs by the BLInG and IG

Network	Objective Value		Test Error	
	BLInG	IG	BLInG	IG
[1 x 50]				
Ailerons	$2,01 \times 10^{-3}$	$1,99 \times 10^{-3}$	$1,98 \times 10^{-3}$	$1,93 \times 10^{-3}$
Bank Marketing	$2,45 \times 10^{-6}$	$2,01 \times 10^{-6}$	$2,83 \times 10^{-5}$	$1,29 \times 10^{-5}$
Bejing PM 2.5	$4,26 \times 10^{-3}$	$4,63 \times 10^{-3}$	$3,98 \times 10^{-3}$	$4,28 \times 10^{-3}$
Bikes Sharing	$1,79 \times 10^{-3}$	$2,21 \times 10^{-3}$	$1,95 \times 10^{-3}$	$2,38 \times 10^{-3}$
California	$1,69 \times 10^{-2}$	$1,89 \times 10^{-2}$	$1,74 \times 10^{-2}$	$1,93 \times 10^{-2}$
CCPP	$3,17 \times 10^{-3}$	$3,21 \times 10^{-3}$	$3,14 \times 10^{-3}$	$3,18 \times 10^{-3}$
Mv	$7,40 \times 10^{-5}$	$8,03 \times 10^{-5}$	$7,72 \times 10^{-5}$	$8,22 \times 10^{-5}$
[3 x 20]				
Ailerons	$2,02 \times 10^{-3}$	$2,05 \times 10^{-3}$	$1,94 \times 10^{-3}$	$2,00 \times 10^{-3}$
Bank Marketing	$2,40 \times 10^{-7}$	$1,26 \times 10^{-6}$	$1,07 \times 10^{-6}$	$4,24 \times 10^{-6}$
Bejing PM 2.5	$4,12 \times 10^{-3}$	$4,19 \times 10^{-3}$	$3,83 \times 10^{-3}$	$3,94 \times 10^{-3}$
Bikes Sharing	$1,78 \times 10^{-3}$	$1,80 \times 10^{-3}$	$1,87 \times 10^{-3}$	$1,93 \times 10^{-3}$
California	$1,61 \times 10^{-2}$	$1,65 \times 10^{-2}$	$1,68 \times 10^{-2}$	$1,74 \times 10^{-2}$
CCPP	$3,14 \times 10^{-3}$	$3,13 \times 10^{-3}$	$3,07 \times 10^{-3}$	$3,07 \times 10^{-3}$
Mv	$1,46 \times 10^{-5}$	$2,10 \times 10^{-5}$	$1,51 \times 10^{-5}$	$2,15 \times 10^{-5}$
[200,50,200]				
Ailerons	$2,09 \times 10^{-3}$	$2,14 \times 10^{-3}$	$2,03 \times 10^{-3}$	$2,07 \times 10^{-3}$
Bank Marketing	$7,14 \times 10^{-6}$	$2,16 \times 10^{-4}$	$1,59 \times 10^{-5}$	$3,22 \times 10^{-4}$
Bejing PM 2.5	$4,60 \times 10^{-3}$	$5,02 \times 10^{-3}$	$4,27 \times 10^{-3}$	$4,68 \times 10^{-3}$
Bikes Sharing	$2,52 \times 10^{-3}$	$3,10 \times 10^{-3}$	$2,53 \times 10^{-3}$	$3,11 \times 10^{-3}$
California	$1,81 \times 10^{-2}$	$2,15 \times 10^{-2}$	$1,90 \times 10^{-2}$	$2,22 \times 10^{-2}$
CCPP	$3,16 \times 10^{-3}$	$3,22 \times 10^{-3}$	$3,10 \times 10^{-3}$	$3,17 \times 10^{-3}$
Mv	$4,87 \times 10^{-5}$	$5,40 \times 10^{-4}$	$4,90 \times 10^{-5}$	$5,33 \times 10^{-4}$
[5 x 50]				
Ailerons	$2,04 \times 10^{-3}$	$2,10 \times 10^{-3}$	$1,99 \times 10^{-3}$	$2,04 \times 10^{-3}$
Bank Marketing	$6,99 \times 10^{-7}$	$1,18 \times 10^{-5}$	$2,46 \times 10^{-6}$	$2,13 \times 10^{-5}$
Bejing PM 2.5	$4,22 \times 10^{-3}$	$4,35 \times 10^{-3}$	$3,92 \times 10^{-3}$	$4,03 \times 10^{-3}$
Bikes Sharing	$2,02 \times 10^{-3}$	$2,19 \times 10^{-3}$	$2,15 \times 10^{-3}$	$2,32 \times 10^{-3}$
California	$1,68 \times 10^{-2}$	$1,77 \times 10^{-2}$	$1,73 \times 10^{-2}$	$1,83 \times 10^{-2}$
CCPP	$3,10 \times 10^{-3}$	$3,15 \times 10^{-3}$	$3,04 \times 10^{-3}$	$3,09 \times 10^{-3}$
Mv	$1,14 \times 10^{-5}$	$5,01 \times 10^{-5}$	$1,14 \times 10^{-5}$	$5,07 \times 10^{-5}$
[10 x 50]				
Ailerons	$2,13 \times 10^{-3}$	$2,17 \times 10^{-3}$	$2,12 \times 10^{-3}$	$2,12 \times 10^{-3}$
Bank Marketing	$8,32 \times 10^{-7}$	$5,01 \times 10^{-5}$	$7,45 \times 10^{-7}$	$4,19 \times 10^{-5}$
Bejing PM 2.5	$4,75 \times 10^{-3}$	$8,68 \times 10^{-3}$	$4,45 \times 10^{-3}$	$8,15 \times 10^{-3}$
Bikes Sharing	$2,88 \times 10^{-3}$	$4,16 \times 10^{-3}$	$2,95 \times 10^{-3}$	$4,31 \times 10^{-3}$
California	$1,80 \times 10^{-2}$	$5,49 \times 10^{-2}$	$1,88 \times 10^{-2}$	$5,81 \times 10^{-2}$
CCPP	$3,30 \times 10^{-3}$	$3,49 \times 10^{-3}$	$3,21 \times 10^{-3}$	$3,40 \times 10^{-3}$
Mv	$3,58 \times 10^{-4}$	$8,31 \times 10^{-4}$	$3,51 \times 10^{-4}$	$7,96 \times 10^{-4}$

6 Conclusions

In this work, we focused on the application of batch and online Block Coordinate Decomposition methods for training Deep Feedforward Neural Networks. We studied how the layered structure of a DFNN can be effectively leveraged for training these models and we defined general *batch* and *minibatch* block layer decomposition schemes. Extensive numerical experiments over different network architectures

Table 7 Cumulative comparison [#wins ; #defs] of BLInG and IG algorithms on training and test error value over the 10 random runs

	Objective value [#wins ; #defs]				
	[1 x 50]	[3 x 20]	[200, 50, 200]	[5 x 50]	[10 x 50]
Ailerons	[0 ; 1]	[0 ; 0]	[2 ; 3]	[1 ; 4]	[9 ; 1]
Bank Marketing	[3 ; 7]	[4 ; 6]	[7 ; 3]	[9 ; 1]	[10 ; 0]
Bejing PM 2.5	[10 ; 0]	[0 ; 1]	[4 ; 5]	[1 ; 2]	[6 ; 0]
Bikes Sharing	[10 ; 0]	[0 ; 4]	[8 ; 2]	[4 ; 2]	[10 ; 0]
California	[8 ; 0]	[0 ; 0]	[7 ; 3]	[1 ; 1]	[10 ; 0]
CCPP	[0 ; 0]	[0 ; 1]	[3 ; 6]	[0 ; 2]	[5 ; 0]
Mv	[5 ; 1]	[5 ; 5]	[9 ; 0]	[10 ; 0]	[10 ; 0]
	Test [#wins ; #defs]				
	[1 x 50]	[3 x 20]	[200, 50, 200]	[5 x 50]	[10 x 50]
Ailerons	[0 ; 3]	[0 ; 0]	[2 ; 3]	[1 ; 4]	[9 ; 1]
Bank Marketing	[3 ; 7]	[3 ; 7]	[8 ; 2]	[10 ; 0]	[10 ; 0]
Bejing PM 2.5	[10 ; 0]	[0 ; 0]	[4 ; 5]	[1 ; 2]	[6 ; 0]
Bikes Sharing	[10 ; 0]	[0 ; 3]	[8 ; 2]	[3 ; 2]	[10 ; 0]
California	[6 ; 0]	[0 ; 0]	[7 ; 3]	[1 ; 1]	[10 ; 0]
CCPP	[0 ; 0]	[0 ; 1]	[4 ; 6]	[2 ; 2]	[4 ; 0]
Mv	[5 ; 1]	[5 ; 5]	[9 ; 0]	[10 ; 0]	[10 ; 0]

Table 8 Ratio between best value found in [10x50] and [1x50]

Network	Objective value		Test Error		
	[1x50]/[10x50]	BLInG	IG	BLInG	IG
Ailerons		1,06	1,09	1,07	1,10
Bank Marketing		0,34	24,91	0,03	3,25
Bejing PM 2.5		1,12	1,88	1,12	1,91
Bikes Sharing		1,61	1,89	1,51	1,81
California		1,06	2,91	1,08	3,01
CCPP		1,04	1,09	1,02	1,07
Mv		4,84	10,35	4,54	9,68

have been performed to assess how performance of state-of-the-art algorithms can be improved. Overall, the application of BCD methods turned out to be effective in avoiding bad attraction regions and speeding up the training process of DFNNs. Both the two proposed methods outperformed no-variable-decomposed counterparts leading to better solutions with better generalization properties as well.

Appendix

Before providing convergence proof of algorithm BLD, we need to recall some properties guaranteed by Armijo Linesearch whose proof can be found here [13].

Lemma 1 *Armijo Linesearch determines a stepsize $\alpha^k \in (0, a]$ within a finite number of iterations. Furthermore, given a sequence $\{w^k\}$ with an accumulation point \bar{w} , if the following limit holds*

$$\lim_{k \rightarrow \infty} f(w^k) - f(w_1^k, \dots, w_\ell^k + \alpha^k d_\ell^k, \dots, w_L) = 0$$

then $\nabla_\ell f(\bar{w}) = 0$

First we show that an Armijo Linesearch coupled with a steepest descent method is always able to find a point \tilde{w} such that conditions 4 and 5 are satisfied.

Theorem 2 *Given a function $f : \mathbb{R}^n \rightarrow [0, \infty)$, conditions (4) and (5) can be satisfied by the application of a finite number n of iterations of the gradient method with an Armijo Linesearch.*

Proof (4) is trivially satisfied $\forall k \geq 1$ since $\nabla f(w)^T d^k = -\|\nabla f(w)\|^2 \leq 0 \forall k$. Concerning condition (5), since α^k is chosen according to Armijo Linesearch, we have that

$$f(w^{k+1}) \leq f(w^k) - \gamma \alpha^k \|\nabla f(w^k)\|^2 \quad \gamma \in (0, 1)$$

Iterating this condition for n steps we get:

$$f(w^{k+n}) \leq f(w^k) - \gamma \left(\sum_{i=0}^{n-1} \alpha^{k+i} \|\nabla f(w^{k+i})\|^2 \right) \quad (9)$$

Since $\alpha^k \nabla f(w^k) = w^k - w^{k+1}$ we have that

$$\alpha^k \|\nabla f(w^k)\|^2 = \frac{\|w^k - w^{k+1}\|^2}{\alpha^k} \quad (10)$$

using (9) and (10) we obtain:

$$f(w^{k+n}) \leq f(w^k) - \gamma \sum_{i=0}^{n-1} \frac{\|w^{k+i} - w^{k+1+i}\|^2}{\alpha^{k+i}} \quad (11)$$

By triangle inequality we have that:

$$\|w^{k+n} - w^k\|^2 \leq \sum_{i=0}^{n-1} \|w^{k+i} - w^{k+1+i}\|^2$$

namely:

$$-\sum_{i=0}^{n-1} \|w^{k+i} - w^{k+1+i}\|^2 \leq -\|w^{k+n} - w^k\|^2$$

which implies:

$$\begin{aligned} f(w^{k+n}) &\leq f(w^k) - \gamma \frac{1}{\alpha^{k+i}} \|w^{k+n} - w^k\|^2 \\ &\leq f(w^k) - \sigma \left(\|w^{k+n} - w^k\|^2 \right) \end{aligned}$$

The function $\sigma(t^k) = \gamma \frac{1}{\alpha^{k+i}} (t^k)^2$ is a forcing function because:

$$\lim_{k \rightarrow \infty} \gamma \frac{1}{\alpha^{k+i}} (t^k)^2 = 0 \implies \lim_{k \rightarrow \infty} (t^k)^2 = 0$$

indeed $\gamma > 0$ and the stepsize $\alpha^k \in (0, a] \forall k$.

We provide the proof of Theorem 1 for the convergence of the BLD method to stationary non minima points with a *forward* cyclical selection of the blocks. This proof can be easily extended to other selections rules by tedious variations of the one provided below.

Proof of Theorem 1 With an abuse of notation, we will refer to $z_{\ell+1}^k$ as the update at iteration k where all the blocks have been updated until block ℓ .

$$z_{\ell+1}^k = (w_1^{k+1}, \dots, w_\ell^{k+1}, w_{\ell+1}^k, \dots, w_L^k)$$

Note that $z_1^k = w^k$ and $z_{L+1}^k = w^{k+1}$. Thanks to condition 4 we have that

$$f(z_{\ell+1}^k) \leq f(w_1^{k+1}, \dots, w_\ell^{k+1} + \alpha^k d_\ell^k, \dots, w_L^k) \leq f z_\ell^k \quad (12)$$

which implies

$$f(w^{k+1}) \leq f(z_\ell^k) \leq f(w^k). \quad (13)$$

Being $f(w)$ defined as the sum of non negative terms we have that $f(w^k) \geq 0$, and being the function coercive thanks to the regularization term $\rho \|w\|^2$, we have that its level curves are compact. Hence, being the sequence $f(w^k)$ a bounded decreasing sequence over a compact set we have that

$$\lim_{k \rightarrow \infty} f(w^k) = \bar{f}. \quad (14)$$

(14) and (12) imply that also the intermediate updates goes to the same limit

$$\lim_{k \rightarrow \infty} f(z_\ell^k) = \bar{f} \quad \forall \ell = 1, \dots, L \quad (15)$$

As a consequence, by (12) we have that

$$\lim_{k \rightarrow \infty} f(z_{\ell+1}^k) - f(w_1^{k+1}, \dots, w_\ell^{k+1} + \alpha^k d_\ell^k, w_{\ell+1}^k, \dots, w_L^k) = 0 \quad (16)$$

Suppose that for a fixed index $\ell \in \{1, \dots, L\}$ z_ℓ^k has an accumulation point \bar{z}_ℓ , Thanks to Armijo Linesearch we have that

$$\nabla_\ell f(\bar{z}_\ell) = 0 \quad \forall \ell = 1, \dots, L \quad (17)$$

Thanks to condition (5) and (14) we have that the distance between two iterates goes to zero

$$\lim_{k \rightarrow \infty} \|z_{\ell+1}^k - z_\ell^k\| = 0 \quad (18)$$

We can consider an accumulation point of w^k , \bar{w} . Since $w^k = z_1^k$, thanks to (18) we have that \bar{w} is also an accumulation point for all the subsequences $z_\ell^k, \ell \in \{1, \dots, L\}$. Thanks to (17) we can conclude that \bar{w} is a stationary point.

References

1. A. Beck and L. Tetrushvili. On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization*, 23(4):2037–2060, 2013.
2. D. P. Bertsekas. Incremental least squares methods and the extended kalman filter. *SIAM J. on Optimization*, 6(3):807–822, Mar. 1996.
3. D. P. Bertsekas. Incremental gradient, subgradient, and proximal methods for convex optimization: A survey. *CoRR*, abs/1507.01030, 2015.
4. D. P. Bertsekas and J. N. Tsitsiklis. Gradient Convergence in Gradient methods with Errors. *SIAM Journal on Optimization*, 10(3):627–642, 2000.
5. L. Bottou. Large-scale machine learning with stochastic gradient descent. In *COMPSTAT*, 2010.
6. L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018.
7. L. Bravi and M. Sciandrone. An incremental decomposition method for unconstrained optimization. *Applied Mathematics and Computation*, 235:80–86, 2014.
8. C. Buzzi, L. Grippo, and M. Sciandrone. Convergent decomposition techniques for training rbf neural networks. *Neural Computation*, 13(8):1891–1920, 2001.

9. V. K. Chauhan, K. Dahiya, and A. Sharma. Mini-batch block-coordinate based stochastic average adjusted gradient methods to solve big data problems. In *Proceedings of the Ninth Asian Conference on Machine Learning*, volume 77 of *Proceedings of Machine Learning Research*, pages 49–64. PMLR, 15–17 Nov 2017.
10. Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in Neural Information Processing Systems 27*, pages 2933–2941. 2014.
11. A. Defazio, F. Bach, and S. Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems 27*, pages 1646–1654. 2014.
12. J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, July 2011.
13. L. Grippo, A. Manno, and M. Sciandrone. Decomposition Techniques for Multilayer Perceptron Training. *IEEE Transactions on Neural Networks and Learning Systems*, 27(11):2146–2159, 2016.
14. L. Grippo and M. Sciandrone. Globally convergent block-coordinate techniques for unconstrained optimization. *Optimization Methods and Software*, 10(4):587–637, 1999.
15. H. Guang-bin, Z. Qin-yu, and S. Chee-kheong. Extreme learning machine: Theory and applications, 2006.
16. G.-B. Huang, D. H. Wang, and Y. Lan. Extreme learning machines: a survey. *International Journal of Machine Learning and Cybernetics*, 2(2):107–122, Jun 2011.
17. R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems 26*, pages 315–323. 2013.
18. D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
19. Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
20. J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer, New York, NY, USA, second edition, 2006.
21. L. Palagi. Global optimization issues in deep network regression: an overview. *Journal of Global Optimization*, pages 1–39, 2018.
22. T. Qin, K. Scheinberg, and D. Goldfarb. Efficient block-coordinate descent algorithms for the group lasso. 5, 06 2013.
23. H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.
24. H. Wang and A. Banerjee. Randomized block coordinate descent for online and stochastic optimization. *arXiv preprint arXiv:1407.0107*, 2014.
25. S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
26. A. W. Yu, L. Huang, Q. Lin, R. Salakhutdinov, and J. Carbonell. Normalized gradient with adaptive stepsize method for deep neural network training. *CoRR*, abs/1707.04822, 2017.
27. T. Zhao, M. Yu, Y. Wang, R. Arora, and H. Liu. Accelerated mini-batch randomized block coordinate descent method. In *Advances in neural information processing systems*, pages 3329–3337, 2014.