Computational Mathematics for Learning and Data Analysis

Implementation of a Neural Network optimized through Stochastic Gradient Descent and Conjugate Gradient Descent



Sabrina Briganti - 465214 - sabrinabriganti@gmail.com Gianmarco Ricciarelli - 555396 - gianmarcoricciarelli@gmail.com

Contents

C	onten	its																			1
1	The 1.1 1.2 1.3 1.4	The no The ba	etwo etwo ack-	ork' pro	s in	nitia gatio	liza on a	tio algo	n . oritl	 hm								 			2
2	Opt 2.1	Stocha 2.1.1 2.1.2 Nonlin 2.2.1 2.2.2 2.2.3	Astic Mo Re near Sea Be	ome gol Co arcl ta	entu ariz onju n Di	ım zzazi	 ione e Gr tion	 e . rad n .	ien	 t . 	· · · · · · · · · · · · · · · · · · ·	 	 	 	 	 		 	 	 	 4 4 5 5
3	Test 3.1 3.2	Monk Cup																			
Bi	bliog	raphy	,																		11

Chapter 1

The network

In this first chapter, we provide some informations about the Artificial Neural Network, i.e. a fully connected Multilayer Perceptron, we implemented from scratch. We'll describe both the network's structure and the algorithm we used in order to make our network learn from the data used during the testing and validation phases. Finally we'll present the loss function we have chosen for our network, and we'll provide and explanation on how it is differentiable. We'll use the notation proposed in [2].

1.1 The network's structure

Since we have to write from scratch an Artificial Neural Network, ANN for short, we have considered some alternatives before choosing the network's final structure. We agreed on a structure composed by one input layer, two hidden layers and one output layer. As convention, the number of units in the input layer is egual to the number of features of the dataset that is used for the learning, validation and testing phases. The two hidden layers contain, respectively, four and eight hidden neurons, following the convention of putting an increasing series of powers of two as number of hidden units per layer. The number of neurons for the output layer depends on the kind of task the network is trying to fullfil. In the case of a classification task, like the MONKS dataset [1], we have decided to put one unit in the output layer, while in the case of a regression task, like the CUP dataset, we have decided to put two units in the output layer. As we have seen studying the papers and books for gathering the necessary knowledge for the project, as [2, 3, 5], choosing to consider the network's structure as an hyperparameter, that is, a variable, could lead to a series of difficult choices during the validation phase, so we have decided to fix the ANN structure to the one described for both the task we have to fullfil, changing only the number of units in the output layer from task to task.

1.2 The network's initialization

1.3 The back-propagation algorithm

The learning procedure for our ANN essentialy consist in two distinct phases:

- 1. compute the network's *gradient*, that is, the derivative of the cost function $\nabla_{\theta} J(\theta)$ with respect to every network's unit;
- 2. optimize the information gathered during the first phase using a distinct technique, like the ones described in chapter 2;

For computing the gradient we have chosen to use the well known backpropagation algorithm, firstly introduced in [6] and described in [2, 3, 5]. This algorithm is also composed by two phases, a first phase, that is, the forward propagation, in which the feature vector \mathbf{x} given in input has to flow from the input layer through the hidden layers and, finally, the output layer, giving the approximation $\hat{\mathbf{y}}$ as output, and a second one, that is, the back-propagation, which allows the informtion to flow backward through the network in order to compute the gradient by applying the Chain Rule of Calculus, that is, a formula for computing the derivative of a composition of functions. It is import to note that with the term back-propagation we mean only the method for computing the gradient, not the whole learning algorithm. We now provide the pseudocode for the forward propagation and the back-propagation phases.

Algorithm 1 Forward propagation through a typical (deep) neural network and the computation of the cost function.

```
1: procedure Forward propagation(l, \mathbf{W}^{(i)} i \in \{1, ..., l\}, \mathbf{b}^{(i)} i \in \{1, ..., l\}, \mathbf{x}, \mathbf{y})

2: \mathbf{h}^{(0)} = \mathbf{x}

3: for k = 1, ..., l do

4: \mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{k-1}

5: \mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})

6: \hat{\mathbf{y}} = \mathbf{h}^{(l)}

7: J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)
```

Algorithm 2 Backward computation for the (deep) neural network of algorithm 1.

```
1: procedure Backward propagation

2: \mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})

3: for k = l, l - 1, \dots, 1 do

4: \mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})

5: \nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)

6: \nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)T} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)

7: \mathbf{g} = \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)T} \mathbf{g}
```

1.4 Loss function is differentiable?

Chapter 2

Optimizers

2.1 Stochastic Gradient Descent

When choosing an optimizer, the *Stochastic Gradient Descent*, SGD for short, is a quite common choice. It is not the best though, since, as proved by the last developments in the machine learning field, its convergence's rate is quite slow. We have implemented a standard SGD version, as described in [2], supporting both *momentum* and *regularization*. [4]

2.1.1 Momentum

2.1.2 Regolarizzazione

2.2 Nonlinear Conjugate Gradient

An intresting optimization ables to lead to an improvement of the performances of the Neural Network, is the use of high-order information during the training phase.

We can, indeed, approximate the *loss function* in a quadratic form, around a given point \mathbf{W} , using the Taylor approximation:

$$\mathcal{E}(\mathbf{W} + \Delta \mathbf{W}) = \mathcal{E}(\mathbf{W}) + \mathbf{g}\mathcal{E}^T \Delta \mathbf{W} + \frac{1}{2} \Delta \mathbf{W}^T \mathbf{H} \Delta \mathbf{W}, \tag{2.1}$$

where \mathbf{H} is the *Hessian* matrix and \mathbf{g} the gradient vector, getting the benefit of choosing the search direction and the step size more carefully by using information from the second order approximation.

Given Eq. 2.1, the optimum adjustment to apply to the weights of the network should be

$$\Delta \mathbf{W}^* = \mathbf{H}^{-1} \mathbf{g},\tag{2.2}$$

carrying the weight of the computation of the inverse of the Hessian.

In order to avoid this expensive computation, we can use the *Conjugate Gradient* methods, which are a class of iterative second-order optimization methods, derived from the steepest-descent algorithm, that ensure low memory requirements.

In this way, the adjustment to the synaptic weights of the network is computed as:

$$\Delta \mathbf{W} = \alpha \mathbf{d},\tag{2.3}$$

where α is the learning rate and **d** is the new direction found.

In our case, the nonlinear conjugate gradient methods are designed to solve the following minimization problem:

$$\min_{\mathbf{W} \in \mathbb{R}^n} \quad \mathcal{E}(\mathbf{W}), \tag{2.4}$$

where \mathcal{E} is the loss function and \mathbf{W} are the synaptic weights of the network.

As showed in the pseudocode ??, the iterative formula generates a sequence of weights $\{W_k\}$, for every epoch of training k, as:

$$\mathbf{W}_{k+1} = \mathbf{W}_k + \alpha \mathbf{d}_k, \quad k = 0, 1, ...,$$
 (2.5)

where α_k is a learning rate and \mathbf{d}_k is a descent direction. These are the new synaptic weights computed with the adjustment of eq.2.3.

2.2.1 Search Direction

The direction \mathbf{d}_k holds the sequent property:

$$\mathbf{d}_k^T \mathbf{H} \mathbf{d}_{tk-1} = 0, \tag{2.6}$$

that means it is conjugate to the previous direction \mathbf{d}_{k-1} . Furthermore, it doesn't need to know all the previous directions, but it only needs the last one, which is why it requires very little storage and computation.

When dealing with quadratic functions, this method keeps the progress obtained so far in the minimization of the loss function, by ensuring that the gradient along the previous direction does not increase. Anyway, it's worth to underline that this method can also be applyed with nonlinear functions: in this case, it should be necessary to restart the process, since there is no assumption that the conjugate directions previously found are still at the minimum of the function.

Each new direction it's a linear combination of the steepest descent -**g** and the previous direction \mathbf{d}_{k-1} , and it is defined as:

$$\mathbf{d}_k = \begin{cases} -\mathbf{g}_0, & \text{if } k = 0; \\ -\mathbf{g}_k + \beta_k \mathbf{d}_{k-1}, & \text{otherwise,} \end{cases}$$
 (2.7)

where β_k is a scalar, to be determined, that says how much of the previous direction should be added to the newest one. When applied to minimize a strictly convex quadratic function, it ensure that the directions \mathbf{d}_k and \mathbf{d}_{k-1} are conjugate with respective to the Hessian of the objective function, that is the property 2.2.1 holds.

Of course, the first search direction when k = 0 is defined as the steepest descent direction at the initial weight \mathbf{W}_0 while, for k > 1, a minimization along each of the search direction is performed.

Since it may be that the direction found is not a descent direction of the objective function, another modified search direction, proposed by Zang et al.??, has been tested in the project. It ensures sufficient descent $g_k^T = -\|g_k\|^2$, indipendent of the line search used or the convexity of the objective function, and is defined as follows:

$$\mathbf{d}_{k}^{+} = \begin{cases} -\mathbf{g}_{0}, & \text{if } k = 0; \\ -(1 + \beta_{k} \frac{\mathbf{g}_{k}^{T} \mathbf{d}_{k}}{\|\mathbf{g}_{k}\|}) \mathbf{g}_{k} + \beta_{k} \mathbf{d}_{k-1}^{+}, & \text{otherwise.} \end{cases}$$
(2.8)

2.2.2 Beta

What really makes the difference in the computation of the conjugate gradient algorithm, is the choice of the method used to compute the β coefficient.

Infact, there has been proposed various choices for computing it, each one giving different efficiency and properties.

The formulas tested in our implementation are four: the Fletcher-Reeves (FR), the Polak-Ribierère (PR), the Hestenes-Stiefel (HS) and a Modified Hestenes-Stiefel (MHS^+).

One of the properties that must be garanteed, is the global convergence of the method. Since, in our network, we are dealing with a nonquadratic loss function, the direction computed as in eq. 2.7 could not be a descent direction. In order to avoid this issue, all the methods have been modified as follows, ensuring the global convergence:

$$\beta^+ = \max\{\beta, 0\}. \tag{2.9}$$

This change provides a sort of restart of the algorithm, in case the β found is negative. This is equivalent to forget the last search direction and start again the search from the steepest descent direction. Furthermore, because of the nonquadratic nature of the error function, the algorithm will not necessarily converge in N steps, as it usually does when applied to quadratic functions. The use of β in eq. 2.9 is similar to adopt the strategy of restarting the algorithm after N steps, initializing d_k to the current steepest descent direction.

$$\beta_k^{PR} = \frac{\mathbf{g}_k^T(\mathbf{g}_k - \mathbf{g}_{k-1})}{\|\mathbf{g}_{k-1}\|^2}, \ \beta_k^{FR} = \frac{\|\mathbf{g}_k\|^2}{\|\mathbf{g}_{k-1}\|^2}, \ \beta_k^{HS} = \frac{\mathbf{g}_k^T(\mathbf{g}_k - \mathbf{g}_{k-1})}{(\mathbf{g}_k - \mathbf{g}_{k-1}^T \mathbf{d}_{k-1})}. \tag{2.10}$$

The HS and the PR methods in eq. 2.10 have very similar performances and they are two of the most efficient conjugate gradient methods, but they are not globally convergent for nonlineat function. That's why the modification of eq. 2.9 has been adopted. Moreover, the HS method is considered superior to other methods when applied to nonquadratic functions.

For what concernes the FR method (also described in eq.2.10), it requires a constrain on the parameters of the inexact line search procedure of section 2.2.3, used to identify the right step length α . In particular, it requires that $\sigma_1 < \sigma_2 < 0.5$ in order to garantee that the Armijo Wolfe conditions are satisfied, and it seems to be less efficient and robust than the other methods. Anyway, by imposing this condition, the FR method is globally convergent even when dealing with nonlinear functions.

The last method tested is the MHS^+ (eq. 2.11, a modified version of the Hestenes-Stiefel one. It garantees sufficient descent with inexact line search and is based on a modified secant equation which approximates the second order information of the loss function with high order accuracy. Moreover, it is globally convergent.

It is defined as follows:

$$\beta_k^{MHS} = \frac{\mathbf{g}_k^T \widetilde{y}_{k-1}^*}{\mathbf{d}_{k-1}^T \widetilde{y}_{k-1}^*}.$$
 (2.11)

In order to better understand the formula 2.11, it's important to describe all the components involved in its definition.

When dealing with quasi-Newton methods, an approximation \mathbf{B}_{k-1} of the Hessian of the loss function $\nabla^2 \mathcal{E}_{k-1}$ is update such that \mathbf{B}_k satisfies the secant condition:

$$\mathbf{B}_k(\mathbf{W}_k - \mathbf{W}_{k-1}) = \mathbf{y}_{k-1},\tag{2.12}$$

where \mathbf{y}_{k-1} is defined as $\mathbf{g}_k - \mathbf{g}_{k-1}$. Wei et al. ?? derived a class of modified secant condition:

$$\mathbf{B}_{k-1}(\mathbf{W}_k - \mathbf{W}_{k-1}) = \widetilde{y}_{k-1},\tag{2.13}$$

$$\widetilde{y}_{k-1} = y_{k-1} + \frac{\theta_{k-1}}{(\mathbf{W}_k - \mathbf{W}_{k-1})^T \mathbf{u}} \mathbf{u}, \tag{2.14}$$

with \mathbf{u} a vector satisfying $(\mathbf{W}_k - \mathbf{W}_{k-1})^T \mathbf{u} \neq 0$ and θ_{k-1} defined as:

$$\theta_{k-1} = 2(\mathcal{E}_{k-1} - \mathcal{E}_k) + (\mathbf{g}_k + \mathbf{g}_{k-1})^T (\mathbf{W}_k - \mathbf{W}_{k-1}). \tag{2.15}$$

Since for $\|(\mathbf{W}_k - \mathbf{W}_{k-1})\| > 1$ the standard secant Eq.2.12 better approximates $\nabla^2 \mathcal{E}_{k-1}(\mathbf{W}_k - \mathbf{W}_{k-1})$ than the modified version in Eq.2.14, Livieris et al. proposed a modification of the equation in this way:

$$\mathbf{B}_{k-1}(\mathbf{W}_k - \mathbf{W}_{k-1}) = \widetilde{y}_{k-1}^*, \tag{2.16}$$

$$\widetilde{y}_{k-1}^* = y_{k-1} + \rho_{k-1} \frac{\max\{\theta_{k-1}, 0\}}{(\mathbf{W}_k - \mathbf{W}_{k-1})^T \mathbf{u}} \mathbf{u}, \tag{2.17}$$

where $\rho_{k-1} \in \{0,1\}$ is a parameter that switch between the standard secant Eq.2.12 and the modified one 2.16, setting ρ_{k-1} as:

$$\rho_{k-1} = \begin{cases} 1, & \|(\mathbf{W}_k - \mathbf{W}_{k-1})\| \le 1; \\ 0, & \text{otherwise.} \end{cases}$$
 (2.18)

It's suggested to use MHS^+ with the search direction \mathbf{d}^+ defined by Eq.2.8.

2.2.3 Line Search

Once computed the new direction **d** involved in the new weights $\mathbf{W} + \alpha \mathbf{d}$, a line search has to be implemented in order to find the right step size which minimize the loss function.

The step size α is nothing more than a scalar: the learning rate for the conjugate gradient algorithm, which tells how far is right to move along a given direction.

So, fixed the values of the weights **W** and the descent direction **d**, the main goal is to find the right value for α that is able to minimize the loss function:

$$\min_{\alpha} \quad \mathcal{E}(\mathbf{W} + \alpha \mathbf{d}). \tag{2.19}$$

Of course, we have to deal with a tradeoff: we want a good reduction, but we can't spend too much time computing the exact value for the optimum solution. So, the smarter way to get it is to use an inexact line search, that try some candidate step size and accepts the first one satisfying some conditions.

This search is performed in two phases:

- a bracketing phase, that finds an initial interval containing a minimizer;
- an interpolation phase that, given the interval, finds the right step length in it.

We decided to use one of the most popular line search condition: the *Armijo-Wolfe* condition.

The search for the better α is led by two condition:

• the Armijo one:

$$\mathcal{E}(W_k + \alpha_k d_k) \le \mathcal{E}(W_k) + \sigma_1 \alpha \nabla \mathcal{E}_k^T d_k \tag{2.20}$$

which ensure that α gives a sufficient decrease of the objective function, being this reduction proportional to the step length α and the directional derivative $\nabla \mathcal{E}_k^T d_k$.

The constant σ_1 has been set $\sigma_1 = 10^{-4}$, since it is suggested in literature to be quite small.

• the Strong Wolfe condition:

$$|\nabla \mathcal{E}(W_k + \alpha_k d_k)^T d_k| \le \mathcal{E}(W_k) + \sigma_2 |\nabla \mathcal{E}_k^T d_k|$$
(2.21)

which garantees to choose steps whose size is not too small.

It is also known as curvature condition and ensures that, moving of a step α along the given direction, the slope of our function if greater than σ_2 times the original gradient (if the slope is only slightly negative, the function cannot decrease rapidly along that direction, so it's better to stop the search).

In this case, the constant σ_2 is equal to 0.1, since a smaller value gives a more accurate line search. Futhermore, having choosen the strong condition, which doesn't allow the derivative to be too positive, we are sure that the α found lies close to a stationary point of the function.

The algorithm satisfing the Strong Wolfe conditions is implemented through three functions, as described in the pseudocodes 3, ??, 5: line_search, zoom and interpolate_alpha. Since two consecutive values may be similar in finite-precision arithmetic, we set a threshold in both the line_search and interpolate_alpha functions, which garantees that the algorithm stops if two values of α are too close or if the maximum number of iterations has been reached.

The line_search function try to find and return a good α ; if it fails, it returns an interval in which continue the searching, invoking the zoom function, which decreases the size of the interval, until it finds and returns a good step length.

Zoom invokes another function, interpolate_alpha, which is nothing more than the implementation of a bisection interpolation in order to find a trial α inside the given interval.

Algorithm 3 Line Search

```
1: procedure LINE SEARCH
 2:
  3:
              \alpha_0 \leftarrow \theta;
  4:
              i \leftarrow 1;
              while i \leq max iter do
  5:
                     Evaluate \mathcal{E}(\alpha_i);
  6:
                     if [\mathcal{E}(\alpha_i) > \mathcal{E}(0) + \sigma_1 \alpha_i \nabla \mathcal{E}_0^T d_0] or [\mathcal{E}(\alpha_i) \leq \mathcal{E}(\alpha_{i-1}) and i > 1] then
  7:
 8:
                            \alpha_* \leftarrow \mathbf{zoom}(\alpha_{i-1}, \alpha_i); \mathbf{return} \ \alpha_*;
                     Evaluate \nabla \mathcal{E}_i
 9:
                     if |\nabla \mathcal{E}_i| \leq -\sigma_2 \nabla \mathcal{E}_0^T d_0 then
10:
                            \alpha_* \leftarrow \alpha_i; return \alpha_*;
11:
                     if \nabla \mathcal{E}_i \geq 0 then
12:
                            \alpha_* \leftarrow \mathbf{zoom}(\alpha_i, \alpha_{i-1}); \mathbf{return} \ \alpha_*;
13:
                     if (|\mathcal{E}_i - \mathcal{E}_{i-1}| \leq threshold then
14:
                            \alpha_* \leftarrow \alpha_i \text{ return } \alpha_*;
15:
                     Choose \alpha_{i+1} \in (\alpha_i, \alpha_{max});
16:
17:
                     i \leftarrow i + 1;
```

Algorithm 4 Zoom

```
1: procedure ZOOM
                while True do
  3:
                       \alpha_i \leftarrow interpolate\_alpha(\alpha_{lo}, \alpha_{hi});
                       Evaluate \mathcal{E}(\alpha_i);
  4:
                       if \left[\mathcal{E}(\alpha_j) > \mathcal{E}(0) + \sigma_1 \alpha_j \nabla \mathcal{E}_0^T d_0\right] or \left[\mathcal{E}(\alpha_j) \leq \mathcal{E}(\alpha_{lo})\right] then
  5:
                               \alpha_* \leftarrow \alpha_i;
  7:
                               return \alpha_*;
                       else
  8:
                               Evaluate \nabla \mathcal{E}_i^T d_j;
  9:
                              if \left|\nabla \mathcal{E}_{j}^{T} d_{j}\right| \leq -\sigma_{2} \nabla \mathcal{E}_{0}^{T} d_{0} then \alpha_{*} \leftarrow \alpha_{j};
10:
11:
12:
                                       return \alpha_*;
                               if \nabla \mathcal{E}_{i}^{T} d_{j}(\alpha_{hi} - \alpha_{lo}) \geq 0 then
13:
14:
                                       \alpha_{hi} \leftarrow \alpha_{lo};
                               if (|\mathcal{E}_j - \mathcal{E}_0| \leq threshold then
15:
16:
                                       \alpha_* \leftarrow \alpha_j
17:
                                       return \alpha_*;
18:
                       \alpha_{lo} \leftarrow \alpha_{j};
```

Algorithm 5 Interpolate

```
1: procedure INTERPOLATE_ALPHA
            i \leftarrow 1;
 2:
            while i \leq max\_iter do
 3:
 4:
                 \alpha_{mid} \leftarrow (\alpha_{hi} - \alpha_{lo})/2
                 Evaluate \mathcal{E}(\alpha_{mid});
 5:
                 if [\mathcal{E}(\alpha_{mid}) == 0] or [(\alpha_{hi} - \alpha_{lo})/2 < threshold] then return \alpha_{mid};
 6:
                 Evaluate \mathcal{E}(\underline{midalpha_{lo}});
 7:
                 if sign(\mathcal{E}(\alpha_{mid})) == sign(\mathcal{E}(\alpha_{lo})) then
 8:
 9:
                       \alpha_{lo} \leftarrow \alpha_{mid};
10:
11:
                       \alpha_{hi} \leftarrow \alpha_{mid};
                 i \leftarrow i + 1;
12:
```

Chapter 3

Test

- 3.1 Monk
- 3.2 Cup

Bibliography

- [1] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org.
- [3] S.S. Haykin. Neural Networks and Learning Machines. Pearson International Edition. Pearson, 2009.
- [4] Ioannis E. Livieris and Panagiotis Pintelas. A new conjugate gradient algorithm for training neural networks based on a modified secant equation. *Applied Mathematics and Computation*, 221:491 502, 2013.
- [5] T.M. Mitchell. *Machine Learning*. McGraw-Hill International Editions. McGraw-Hill, 1997.
- [6] D. E. RUMERLHAR. Learning representation by back-propagating errors. *Nature*, 323:533-536, 1986.