

## Gradient methods for parameter optimization

### Exercise T4.1: Multilayer perceptron

(tutorial)

- (a) Recap the optimization of the MLP parameters (via the backpropagation algorithm).
- (b) Outline the weight space symmetries giving rise to  $\prod_{v=1}^L N_v! \cdot 2^{N_v}$  equivalent solutions where  $L$  is the number of hidden layers and  $N_v$  the respective number of neurons in layer  $v \implies$  no unique global minimum but a large equivalence class of (best) solutions.

### Exercise T4.2: Linear neuron for regression

(tutorial)

To prepare for the homework, we discuss a simple connectionist neuron with linear output function for a real one-dimensional input  $x \in \mathbb{R}$  and output  $y \in \mathbb{R}$ .

- (a) Describe the output function  $y(x; \underline{\mathbf{w}})$  of the neuron in vector notation.
- (b) Derive the gradient and Hessian matrix of the quadratic error function.
- (c) Solve the optimization of the quadratic error function for a data set  $\{(x^{(\alpha)}, y_T^{(\alpha)})\}_{\alpha=1, \dots, p}$  analytically in matrix form.
- (d) Calculate the solution when the objective includes the quadratic training cost  $E^T$  plus a “weight decay” regularization term as used in *ridge regression*, i.e.

$$R_{[\underline{\mathbf{w}}]} = E_{[\underline{\mathbf{w}}]}^T + \lambda \|\underline{\mathbf{w}}\|^2$$

### Exercise T4.3: Conjugate gradient

(tutorial)

- (a) How does the convergence speed of *gradient descent* depend on the learning rate  $\eta$ ?
- (b) Describe how *line search* speeds up convergence.
- (c) What is a *conjugate direction* and how can it improve convergence speed?

### Exercise H4.1: Line search

(homework, 4 points)

In this exercise you will analyze line search based on the simple example of a linear neuron with quadratic cost function  $E_{[\underline{\mathbf{w}}]}^T$ . Here we optimize the cost function along a given direction  $\underline{\mathbf{d}}_t$  (that can be but is not necessarily identical to the gradient):

$$\underline{\mathbf{w}}_{t+1} = \underline{\mathbf{w}}_t - \eta_t \underline{\mathbf{d}}_t.$$

- (a) (1 point) Derive the 2<sup>nd</sup> order Taylor approximation of an arbitrary  $E_{[\underline{\mathbf{w}}_{t+1}]}^T$  around  $\underline{\mathbf{w}}_t$ .
- (b) (1 point) Derive a bound on the step size  $\eta_t$  using the above approximation in

$$E_{[\underline{\mathbf{w}}_{t+1}]}^T \stackrel{!}{\leq} E_{[\underline{\mathbf{w}}_t]}^T.$$

- (c) (1 point) Derive the optimal step size  $\eta_t^*$  for the quadratic cost function

$$E_{[\mathbf{w}]}^T = \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

by minimizing the cost function w.r.t.  $\eta$ . Make sure your solution depends only on known quantities like the weight vector  $\mathbf{w}_t$ , the gradient  $\nabla E_{[\mathbf{w}_t]}^T$  and/or the Hessian  $\mathbf{H}$  of  $E_{[\mathbf{w}_t]}^T$ .

- (d) (1 point) Prove that the gradient  $\nabla E_{[\mathbf{w}_{t+1}]}^T$  after one update step with *line search* is orthogonal to the optimized direction  $\mathbf{d}_t$ .

### Solution

- (a) Let the gradient  $\mathbf{g}_t := \nabla E_{[\mathbf{w}_t]}^T$  and the Hessian matrix  $\mathbf{H}_t := \Delta E_{[\mathbf{w}_t]}^T$ , then

$$\begin{aligned} E_{[\mathbf{w}_{t+1}]}^T &\approx E_{[\mathbf{w}_t]}^T + (\mathbf{w}_{t+1} - \mathbf{w}_t)^\top \mathbf{g}_t + \frac{1}{2}(\mathbf{w}_{t+1} - \mathbf{w}_t)^\top \mathbf{H}_t(\mathbf{w}_{t+1} - \mathbf{w}_t) \\ &= E_{[\mathbf{w}_t]}^T - \eta_t \mathbf{d}_t^\top \mathbf{g}_t + \frac{\eta_t^2}{2} \mathbf{d}_t^\top \mathbf{H}_t \mathbf{d}_t. \end{aligned}$$

$$(b) \quad E_{[\mathbf{w}_{t+1}]}^T \approx E_{[\mathbf{w}_t]}^T - \eta_t \mathbf{d}_t^\top \mathbf{g}_t + \frac{\eta_t^2}{2} \mathbf{d}_t^\top \mathbf{H}_t \mathbf{d}_t \stackrel{!}{\leq} E_{[\mathbf{w}_t]}^T \quad \Rightarrow \quad \eta_t \stackrel{!}{\leq} 2 \frac{\mathbf{d}_t^\top \mathbf{g}_t}{\mathbf{d}_t^\top \mathbf{H}_t \mathbf{d}_t}.$$

- (c) Solve  $\min_{\eta} E_{[\mathbf{w}_{t+1}]}^T$  by setting the derivative w.r.t.  $\eta$  to zero:

$$\begin{aligned} \frac{\partial E_{[\mathbf{w}_{t+1}]}^T}{\partial \eta} &= \left( \frac{\partial E_{[\mathbf{w}_{t+1}]}^T}{\partial \mathbf{w}_{t+1}} \right)^\top \frac{\partial \mathbf{w}_{t+1}}{\partial \eta} = (\mathbf{H}_t \mathbf{w}_t - \eta_t \mathbf{H}_t \mathbf{d}_t - \mathbf{H}_t \mathbf{w}^*)^\top (-\mathbf{d}_t) \stackrel{!}{=} 0 \\ \Rightarrow \quad \eta^* &= \frac{\mathbf{d}_t^\top \mathbf{H}_t (\mathbf{w}_t - \mathbf{w}^*)}{\mathbf{d}_t^\top \mathbf{H}_t \mathbf{d}_t} = \frac{\mathbf{d}_t^\top \mathbf{g}_t}{\mathbf{d}_t^\top \mathbf{H}_t \mathbf{d}_t}, \quad \text{as } \mathbf{g}_t = \mathbf{H}_t (\mathbf{w}_t - \mathbf{w}^*). \end{aligned}$$

- (d) The gradient is orthogonal to the direction if  $\mathbf{d}_t^\top \mathbf{g}_{t+1} = 0$ .

$$\begin{aligned} \mathbf{g}_{t+1} &= \mathbf{H}_t (\mathbf{w}_{t+1} - \mathbf{w}^*) = \mathbf{H}_t (\mathbf{w}_t - \eta_t \mathbf{d}_t - \mathbf{w}^*) = \mathbf{g}_t - \eta_t \mathbf{H}_t \mathbf{d}_t \\ \mathbf{d}_t^\top \mathbf{g}_{t+1} &= \mathbf{d}_t^\top \mathbf{g}_t - \frac{\mathbf{d}_t^\top \mathbf{g}_t}{\mathbf{d}_t^\top \mathbf{H}_t \mathbf{d}_t} \mathbf{d}_t^\top \mathbf{H}_t \mathbf{d}_t = 0 \end{aligned}$$

**Exercise H4.2: Comparison of gradient descent methods (homework, 6 points)**

In this exercise we compare the performance of three learning procedures applied to a simple connectionist neuron with a linear output function. All procedures will compute the gradient using the entire training set (batch gradient descent). The procedures are: (i) Gradient (or steepest) descent with constant learning rate, (ii) steepest descent combined with a line search method to determine the learning rate, and (iii) the conjugate gradient method.

**Training Data:** The training data set consists of three points ( $p = 3$ ):

$$\{(x^{(\alpha)}, y_T^{(\alpha)})\} = \{(-1, -0.1), (0.3, 0.5), (2, 0.5)\},$$

i.e. for a given data point, both input and output are scalar values.

**Cost function:** The gradient for the *quadratic error* function is given by

$$\underline{\mathbf{g}} = \frac{\partial E^T}{\partial \underline{\mathbf{w}}} = \underline{\mathbf{H}} \underline{\mathbf{w}} + \underline{\mathbf{b}}, \quad \text{with } \underline{\mathbf{H}} = \underline{\mathbf{X}} \underline{\mathbf{X}}^T \quad \text{and} \quad \underline{\mathbf{b}} = -\underline{\mathbf{X}} \underline{\mathbf{y}}^T,$$

where  $\underline{\mathbf{X}} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(p)} \end{pmatrix} \in \mathbb{R}^{2,p}$  and  $\underline{\mathbf{y}} = (y_T^{(1)}, y_T^{(2)}, \dots, y_T^{(p)}) \in \mathbb{R}^{1,p}$ .

**Initialization:** Use the following initialization for all three (batch) gradient methods:

$$\underline{\mathbf{w}}_1 = (w_0, w_1)_1^T = (-0.45, 0.2)^T$$

- (a) (2 points) *Gradient Descent:* Implement a steepest descent procedure where the weights at iteration  $t + 1$  are calculated using the weights and the gradient at iteration  $t$

$$\underline{\mathbf{w}}_{t+1} = \underline{\mathbf{w}}_t - \eta \underline{\mathbf{g}}_t,$$

with an adequate learning rate  $\eta$  and where  $\underline{\mathbf{g}}_t = \underline{\mathbf{g}}(\underline{\mathbf{w}}_t)$ . Plot the resulting weight vectors from all iterations as a scatter plot ( $w_0$  vs.  $w_1$ ), and in an additional plot ( $w_i$  vs. iterations  $t$ ), to show the development of each parameter during gradient descent.

- (b) (2 points) *Line Search:* Implement a line search procedure

$$\underline{\mathbf{w}}_{t+1} = \underline{\mathbf{w}}_t - \eta \underline{\mathbf{g}}_t, \quad \text{with optimal step size} \quad \eta = \frac{\underline{\mathbf{g}}_t^T \underline{\mathbf{g}}_t}{\underline{\mathbf{g}}_t^T \underline{\mathbf{H}} \underline{\mathbf{g}}_t}.$$

Plot the resulting weight vectors from all iterations as a scatter plot ( $w_0$  vs.  $w_1$ ), and in an additional plot ( $w_i$  vs. iterations  $t$ ), to show the development of the parameters during line search.

- (c) (2 points) *Conjugate Gradient:* Implement a conjugate gradient procedure:

Initialize:  $\underline{\mathbf{w}}_1, \underline{\mathbf{d}}_1 = -\underline{\mathbf{g}}_1$

**while** stopping criterion not satisfied **do**

minimize  $E$  along  $\underline{\mathbf{d}}_t$ :  $\underline{\mathbf{w}}_{t+1} = \underline{\mathbf{w}}_t + \eta_t \underline{\mathbf{d}}_t$  with step size  $\eta_t = -\frac{\underline{\mathbf{d}}_t^T \underline{\mathbf{g}}_t}{\underline{\mathbf{d}}_t^T \underline{\mathbf{H}} \underline{\mathbf{d}}_t}$

calculate new gradient  $\underline{\mathbf{g}}_{t+1} = \underline{\mathbf{H}} \underline{\mathbf{w}}_{t+1} + \underline{\mathbf{b}}$

calculate new conjugate direction  $\underline{\mathbf{d}}_{t+1} = \underline{\mathbf{g}}_{t+1} + \beta_t \underline{\mathbf{d}}_t$  with “momentum”

$$\beta_t = -\frac{\underline{\mathbf{g}}_{t+1}^T \underline{\mathbf{g}}_{t+1}}{\underline{\mathbf{g}}_t^T \underline{\mathbf{g}}_t}. \quad (\text{Fletcher-Reeves form})$$

increase  $t \leftarrow t + 1$

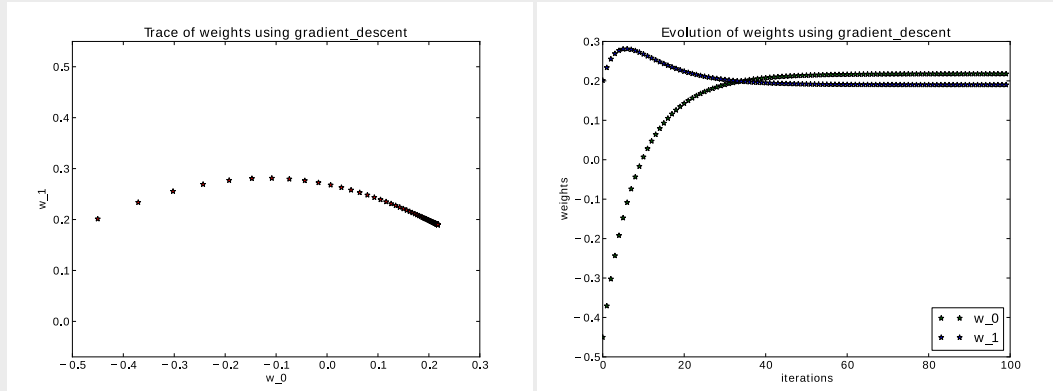
**end**

Plot the resulting weight vectors from all iterations as a scatter plot ( $w_0$  vs.  $w_1$ ), and in an additional plot ( $w_i$  vs. iterations  $t$ ), to show the development of the parameters during conjugate gradient descent.

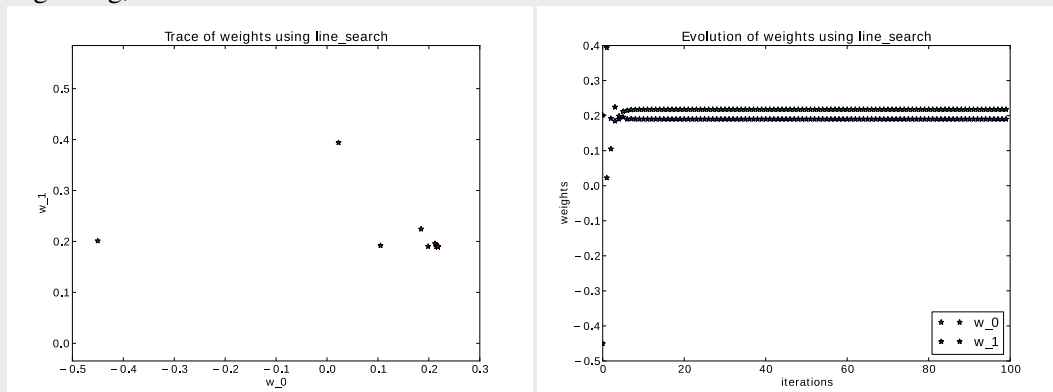
Compare the different methods in terms of convergence behaviour.

### Solution

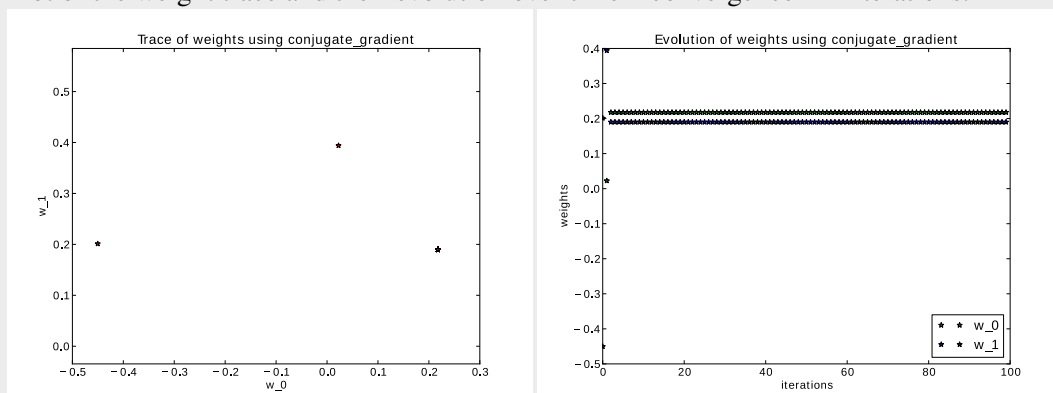
(a) Plot of the weight trace and their evolution over time – fairly slow convergence.



(b) Plot of the weight trace and their evolution over time – much faster convergence in the beginning, noticeable slower close to the minimum.



(c) Plot of the weight trace and their evolution over time – convergence in 2 iterations.



**Total 10 points.**