

## 7 Automatic Differentiation (AutoDiff)

### Goal

Forward and reverse mode auto-differentiation.

### Alert 7.1: Convention

Gray boxes are not required hence can be omitted for unenthusiastic readers.

[This note is likely to be updated again soon.](#)

### Definition 7.2: Function superposition and computational graph (Bauer 1974)

Let  $\mathcal{F}_0$  be a class of *basic* functions. A (vector-valued) function  $g : \mathbb{X} \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^m$  is a **superposition of the basic class  $\mathcal{F}_0$**  if the following is satisfied:

- There exist some DAG  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  where using topological sorting we arrange the nodes as follows:

$$\underbrace{v_1, \dots, v_d}_{\text{input}} \quad \underbrace{v_{d+1}, \dots, v_{d+k}}_{\text{intermediate variables}} \quad \underbrace{v_{d+k+1}, \dots, v_{d+k+m}}_{\text{output}} \quad \text{and } (v_i, v_j) \in \mathcal{E} \implies i < j.$$

Here we implicitly assume the outputs of the function  $g$  do not depend on each other. If they do, we need only specify the indices of the output nodes accordingly (i.e., they may not all appear in the end).

- For each node  $v_i$ , let  $\mathcal{S}_i := \{u \in \mathcal{V} : (u, v_i) \in \mathcal{E}\}$  and  $\mathcal{O}_i := \{u \in \mathcal{V} : (v_i, u) \in \mathcal{E}\}$  denote the (immediate) predecessors and successors of  $v_i$ , respectively. Clearly,  $\mathcal{S}_i = \emptyset$  if  $i \leq d$  (i.e., input nodes) and  $\mathcal{O}_i = \emptyset$  if  $i > d + k$  (i.e., output nodes).
- The nodes are computed as follows: sequentially for  $i = 1, \dots, d + k + m$ ,

$$v_i = \begin{cases} x_i, & i \leq d \\ f_i(\mathcal{S}_i), & i > d \end{cases}, \quad \text{where } f_i \in \mathcal{F}_0. \quad (7.1)$$

Our definition of superposition closely resembles the computational graph of Bauer (1974), who attributed the idea to Kantorovich (1957).

Bauer, F. L. (1974). “Computational Graphs and Rounding Error”. *SIAM Journal on Numerical Analysis*, vol. 11, no. 1, pp. 87–96.

Kantorovich, L. V. (1957). “On a system of mathematical symbols, convenient for electronic computer operations”. *Soviet Mathematics Doklady*, vol. 113, no. 4, pp. 738–741.

### Exercise 7.3: Neural networks as function superposition

Let  $\mathcal{F}_0 = \{+, \times, \sigma, \text{constant}\}$ . Prove that any multi-layer NN is a superposition of the basic class  $\mathcal{F}_0$ .  
Is exp a superposition of the basic class above?

### Theorem 7.4: Automatic differentiation (e.g., Kim et al. 1984)

Let  $\mathcal{F}_0$  be a basic class of *differentiable* functions that includes  $+$ ,  $\times$ , and all constants. Denote  $T(f)$  as the complexity of computing the function  $f$  and  $T(f, \nabla f)$  the complexity with additional computation of the gradient. Let  $\mathcal{I}_f$  and  $\mathcal{O}_f$  be the input and output arguments and assume there exists some constant

$C = C(\mathcal{F}_0) > 0$  so that

$$\forall f \in \mathcal{F}_0, \quad T(f, \nabla f) + |\mathcal{S}_f| |\mathcal{O}_f| [T(+) + T(\times) + T(\text{constant})] \leq C \cdot T(f). \quad (7.2)$$

Then, for any superposition  $g : \mathbb{R}^d \rightarrow \mathbb{R}^m$  of the basic class  $\mathcal{F}_0$ , we have

$$T(g, \nabla g) \leq \begin{cases} C\gamma d \cdot T(g), & \text{forward mode} \\ C\gamma m \cdot T(g), & \text{reverse mode} \end{cases}, \quad \text{i.e., } T(g, \nabla g) \leq C\gamma(m \wedge d) \cdot T(g),$$

where  $\gamma$  is the maximum output dimension of the basic functions used to superpose  $g$ .

*Proof:* Applying the [chain rule](#) to the recursive formula (7.1) it is clear that any superposition  $g$  is differentiable too. We split the proof into two parts: a forward mode and a backward mode.

**Forward mode:** Let us define the block matrix  $U = [U_1, \dots, U_d, U_{d+1}, \dots, U_{d+k}, U_{d+k+1}, \dots, U_{d+k+m}] \in \mathbb{R}^{d \times \sum_i d_i}$ , where each column block  $U_i$  corresponds to the gradient  $\frac{\partial v_i}{\partial \mathbf{x}} \in \mathbb{R}^{d \times d_i}$ , where  $d_i$  is the output dimension of node  $v_i$  (typically 1). By definition of the input nodes we have

$$U_i = \mathbf{e}_i, \quad i = 1, \dots, d,$$

where  $\mathbf{e}_i$  is the standard basis vector in  $\mathbb{R}^d$ . Using the recursive formula (7.1) and chain rule we have

$$U_i = \frac{\partial v_i}{\partial \mathbf{x}} = \sum_{j \in \mathcal{S}_i} \frac{\partial v_j}{\partial \mathbf{x}} \cdot \frac{\partial v_i}{\partial v_j} = \sum_{j \in \mathcal{S}_i} U_j \cdot \nabla_j f_i, \quad \text{where } \nabla_j f_i = \frac{\partial f_i}{\partial v_j} \in \mathbb{R}^{d_j \times d_i}.$$

In essence, [by differentiating at each node, we obtain a square and sparse system of linear equations](#), where  $\nabla_j f_i$  are known coefficients and  $U_i$  are unknown variables. Solving the linear system yields  $U_{d+k+1}, \dots, U_{d+k+m}$ , the desired gradient of  $g$ . Thanks to the topological ordering, we can simply solve  $U_i$  one by one. Let  $\gamma = \max_i d_i$  be the maximum output dimension of any node. We bound the complexity of the forward mode as follows:

$$\begin{aligned} T(g, \nabla g) &\leq \sum_{i \in \mathcal{V}} T(f_i, \nabla f_i) + \sum_{j \in \mathcal{S}_i} dd_i d_j [T(+) + T(\times) + T(\text{constant})] \\ &\leq d\gamma \sum_{i \in \mathcal{V}} T(f_i, \nabla f_i) + |\mathcal{S}_{f_i}| |\mathcal{O}_{f_i}| [T(+) + T(\times) + T(\text{constant})] \leq d\gamma \sum_{i \in \mathcal{V}} C \cdot T(f_i) = d\gamma C \cdot T(g). \end{aligned}$$

**Reverse mode:** Let us rename the outputs  $y_i = v_{d+k+i}$  for  $i = 1, \dots, m$ . Similarly we define the block matrix  $V = [V_1; \dots; V_d; V_{d+1}; \dots; V_{d+k}; V_{d+k+1}; \dots; V_{d+k+m}] \in \mathbb{R}^{\sum_i d_i \times m}$ , where each row block  $V_i$  corresponds to the transpose of the gradient  $\frac{\partial \mathbf{y}}{\partial v_i} \in \mathbb{R}^{m \times d_i}$ , where  $d_i$  is the output dimension of node  $v_i$  (typically 1). By definition of the output nodes we have

$$V_{d+k+i} = \mathbf{e}_i, \quad i = 1, \dots, m, \quad \mathbf{e}_i \in \mathbb{R}^{1 \times m}.$$

Using the recursive formula (7.1) and chain rule we have

$$V_i = \frac{\partial \mathbf{y}}{\partial v_i} = \sum_{j \in \mathcal{O}_i} \frac{\partial v_j}{\partial v_i} \cdot \frac{\partial \mathbf{y}}{\partial v_j} = \sum_{j \in \mathcal{O}_i} \nabla_i f_j \cdot V_j, \quad \text{where } \nabla_i f_j = \frac{\partial f_j}{\partial v_i} \in \mathbb{R}^{d_i \times d_j}. \quad (7.3)$$

Again, [by differentiating at each node we obtain a square and sparse system of linear equations](#), where  $\nabla_i f_j$  are known coefficients and  $V_i$  are unknown variables. Solving the linear system yields  $V_1, \dots, V_d$ , the desired gradient of  $g$ . Thanks to the topological ordering, we can simply solve  $V_i$  one by one [backwards](#), after a forward pass to get the function values at each node. Similar as the forward mode, we can bound the complexity as  $m\gamma C \cdot T(g)$ . ■

Thus, surprisingly, [for real-valued superpositions \( \$m = \gamma = 1\$ \), computing the gradient, which is a  \$d \times 1\$  vector, costs at most constant times that of the function value \(which is a scalar\), if we operate in the reverse mode!](#) The common misconception is that the gradient has size  $d \times 1$  hence if we compute one component

at a time we end up  $d$  times slower. This is wrong, because we can recycle computations. Note also that even reading the input already costs  $O(d)$ . However, this time complexity gain, as compared to that of the forward mode, is achieved through a space complexity tradeoff: in reverse mode we need a forward pass first to collect and store all function values at each node, whereas in the forward mode these function values can be computed on the fly.

We note that [in the proof we \(tacitly\) took for granted that](#)

$$T(g) = \sum_{i \in \mathcal{V}} T(f_i),$$

i.e., the computational cost of the supposition  $g$  is simply the total cost of each component (e.g., by following the recursion (7.1) blindly).

Kim, K. V., Y. E. Nesterov, and B. V. Cherkasskii (1984). “An estimate of the effort in computing the gradient”. *Soviet Mathematics Doklady*, vol. 29, no. 2, pp. 384–387.

### Algorithm 7.5: Automatic differentiation (AD) pseudocode

We summarize the forward and reverse algorithms below. Note that to compute the gradient-vector multiplication  $\mathbf{w}^\top \nabla g$ , we can use the forward mode and initialize  $V_i = w_i$  (i.e., multiplying  $\mathbf{w}$  from left on both sides of line 8). Similarly, to compute  $(\nabla g)\mathbf{w}$ , we can use the reverse mode and initialize  $V_{d+k+i} = w_i$ .

---

#### Algorithm: Forward-mode automatic differentiation for superposition

---

**Input:**  $\mathbf{x} \in \mathbb{R}^d$ , basic function class  $\mathcal{F}_0$ , computational graph  $\mathcal{G}$   
**Output:** gradient  $[V_{d+k+1}, \dots, V_{d+k+m}] \in \mathbb{R}^{d \times m}$

```

1 for  $i = 1, \dots, d$  do // forward: initialize function values and derivatives
2    $v_i \leftarrow x_i$ 
3    $V_i \leftarrow \mathbf{e}_i \in \mathbb{R}^{d \times 1}$ 
4 for  $i = d+1, \dots, d+k+m$  do // forward: accumulate function values and derivatives
5   compute  $v_i \leftarrow f_i(\mathcal{I}_i)$ 
6   for  $j \in \mathcal{I}_i$  do
7     compute partial derivatives  $\nabla_j f_i(\mathcal{I}_i)$ 
8    $V_i \leftarrow \sum_{j \in \mathcal{I}_i} V_j \cdot \nabla_j f_i$ 

```

---



---

#### Algorithm: Reverse-mode automatic differentiation for superposition

---

**Input:**  $\mathbf{x} \in \mathbb{R}^d$ , basic function class  $\mathcal{F}_0$ , computational graph  $\mathcal{G}$   
**Output:** gradient  $[V_1; \dots; V_d] \in \mathbb{R}^{d \times m}$

```

1 for  $i = 1, \dots, d$  do // forward: initialize function values
2    $v_i \leftarrow x_i$ 
3 for  $i = d+1, \dots, d+k+m$  do // forward: accumulate function values
4   compute  $v_i \leftarrow f_i(\mathcal{I}_i)$ 
5 for  $i = 1, \dots, m$  do // backward: initialize derivatives
6    $V_{d+k+i} \leftarrow \mathbf{e}_i \in \mathbb{R}^{1 \times m}$ 
7 for  $i = d+k, \dots, 1$  do // backward: accumulate derivatives
8   for  $j \in \mathcal{O}_i$  do
9     compute partial derivatives  $\nabla_i f_j(\mathcal{I}_i)$ 
10   $V_i \leftarrow \sum_{j \in \mathcal{O}_i} \nabla_i f_j \cdot V_j$ 

```

---

We remark that, as suggested by Wolfe (1982), one effective way to test AD (or manually programmed derivatives) and locate potential errors is through the classic [finite difference](#) approximation.

Wolfe, P. (1982). “Checking the Calculation of Gradients”. *ACM Transactions on Mathematical Software*, vol. 8, no. 4, pp. 337–343.

**Example 7.6: Some applications of AD**

- Directional derivative in line search, such as cubic interpolation, can be cheaply computed using Algorithm 7.5 (Kim et al. 1984).
- Consider a (recurrent) function  $f$  defined through recursion:

$$f(\mathbf{w}) = \phi_t(\mathbf{x}_t, \mathbf{w}), \quad \text{where } \mathbf{x}_{k+1} = \phi_k(\mathbf{x}_k, \mathbf{w}), \quad k = 0, \dots, t-1.$$

The forward-mode differentiation (e.g., Kim et al. 1984, p. 66) is

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = U_t \nabla_{\mathbf{x}} \phi_t(\mathbf{x}_t, \mathbf{w}) + \nabla_{\mathbf{w}} \phi_t(\mathbf{x}_t, \mathbf{w}), \quad \text{where } U_{k+1} := \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{w}} = U_k \nabla_{\mathbf{x}} \phi_k(\mathbf{x}_k, \mathbf{w}) + \nabla_{\mathbf{w}} \phi_k(\mathbf{x}_k, \mathbf{w}),$$

with  $U_0 = I$ , while the reverse-mode differentiation is

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = \sum_{k=0}^t \nabla_{\mathbf{w}} \phi_k(\mathbf{x}_k, \mathbf{w}) V_{k+1}, \quad \text{where } V_k := \frac{\partial f}{\partial \mathbf{x}_k} = \nabla_{\mathbf{x}} \phi_k(\mathbf{x}_k, \mathbf{w}) V_{k+1}, \quad V_{t+1} = I.$$

Kim, K. V., Y. E. Nesterov, V. A. Skokov, and B. V. Cherkasskii (1984). “An efficient algorithm for computing derivatives and extremal problems”. *Ekonomika i matematicheskie metody*, vol. 20, no. 2, pp. 309–318.

**Exercise 7.7: Matrix multiplication**

To understand the difference between forward-mode and backward-mode differentiation, let us consider the simple matrix multiplication problem: Let  $A_\ell \in \mathbb{R}^{d_\ell \times d_{\ell+1}}$ ,  $\ell = 1, \dots, L$ , where  $d_1 = d$  and  $d_{L+1} = m$ . We are interested in computing

$$A = \prod_{\ell=1}^L A_\ell.$$

- What is the complexity if we multiply from left to right (i.e.  $\ell = 1, 2, \dots, L$ )?
- What is the complexity if we multiply from right to left (i.e.  $\ell = L, L-1, \dots, 1$ )?
- What is the optimal way to compute the product?

**Remark 7.8: Further insights on AD**

If we associate an edge weight  $w_{ij} = \frac{\partial v_j}{\partial v_i}$  to  $(i, j) \in \mathcal{E}$ , then the desired gradient

$$\frac{\partial g_i}{\partial x_j} = \sum_{\text{path } P: v_j \rightarrow v_i} \prod_{e \in P} w_e. \quad (7.4)$$

However, we cannot compute the above naively, as the number of paths in a DAG can grow exponentially quickly with the depth. The forward and reverse modes in the proof of Theorem 7.4 correspond to two dynamic programming solutions. (Incidentally, this is exactly how one computes the [graph kernel](#) too.)

Naumann (2008) showed that finding the optimal way to compute (7.4) is NP-hard.

Naumann, U. (2008). “Optimal Jacobian accumulation is NP-complete”. *Mathematical Programming*, vol. 112, no. 2, pp. 427–441.

**Remark 7.9: Tightness of dimension dependence in AD (e.g., Griewank 2012)**

The dimensional dependence  $m \wedge d$  cannot be reduced in general. Indeed, consider the simple function  $\mathbf{f}(\mathbf{x}) = \sin(\langle \mathbf{x}, \mathbf{w} \rangle) \mathbf{b}$ , where  $\mathbf{x} \in \mathbb{R}^d$  and  $\mathbf{b} \in \mathbb{R}^m$ . Computing  $\mathbf{f}$  clearly costs  $O(d + m)$  (assuming  $\sin$  can be evaluated in  $O(1)$ ) while even outputting the gradient costs  $O(dm)$ .

Griewank, A. (2012). “Who Invented the Reverse Mode of Differentiation?” *Documenta Mathematica*, vol. Extra Volume ISMP, pp. 389–400.

**Exercise 7.10: Backpropagation (e.g., Rumelhart et al. 1986)**

Apply Theorem 7.4 to multi-layer NNs and recover the celebrated backpropagation algorithm. Distinguish two cases:

- Fix the network weights  $W_1, \dots, W_L$  and compute the derivative w.r.t. the input  $\mathbf{x}$  of the network. This is useful for constructing [adversarial examples](#).
- Fix the input  $\mathbf{x}$  of the network and compute the derivative w.r.t. the network weights  $W_1, \dots, W_L$ . This is useful for [training](#) the network.

Suppose we know how to compute the derivatives of  $f(x, y)$ . Explain how to compute the derivative of  $f(x, x)$ ?

- Generalize from above to derive the backpropagation rule for convolutional neural nets (CNN).
- Generalize from above to derive the backpropagation rule for recurrent neural nets (RNN).

Rumelhart, D. E., G. E. Hinton, and R. J. Williams (1986). “Learning representations by back-propagating errors”. *Nature*, vol. 323, pp. 533–536.

**Remark 7.11: Fast computation of other derivatives (Kim et al. 1984)**

Kim et al. (1984) pointed out an important observation, namely that the proof of Theorem 7.4 only uses the chain-rule property of differentiation:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial x}.$$

In other words, we **could replace differentiation with any other operation that respects the chain rule and obtain the same efficient procedure for computation**. For instance, the relative differential in numerical analysis or the directional derivative can both be efficiently computed in the same way.

As hinted by Kim et al. (1984, p. 62), generalized derivatives for nonsmooth functions may also be efficiently computed in a similar manner, as long as the chain rule still holds under additional regularity conditions. See Nesterov (1987, 2005) and Kakade and Lee (2018).

Kim, K. V., Y. E. Nesterov, V. A. Skokov, and B. V. Cherkasskii (1984). “An efficient algorithm for computing derivatives and extremal problems”. *Ekonomika i matematicheskie metody*, vol. 20, no. 2, pp. 309–318.

Nesterov, Y. (1987). “The technique of nonsmooth differentiation”. *Engineering Cybernetics: Soviet Journal of Computer and Systems Science*, vol. 25, no. 6, pp. 113–123.

— (2005). “Lexicographic differentiation of nonsmooth functions”. *Mathematical Programming*, vol. 104, pp. 669–700.

Kakade, S. M. and J. D. Lee (2018). “Provably Correct Automatic Sub-Differentiation for Qualified Programs”. In: *NIPS*, pp. 7125–7135.

**Algorithm 7.12: Reverse-mode Hessian-vector product (Kim et al. 1984)**

Similarly, one can compute the Hessian-vector product efficiently as it also respects the chain rule. Indeed, following Kim et al. (1984, p. 61) we denote the directional derivative of  $\mathbf{f}$  along direction  $\mathbf{z}$  as

$$\mathfrak{D}(\mathbf{f}) = \mathfrak{D}(\mathbf{f}(\mathbf{x}); \mathbf{z}) := \left. \frac{d\mathbf{f}(\mathbf{x} + t\mathbf{z})}{dt} \right|_{t=0} = \mathbf{f}'(\mathbf{x})(\mathbf{z}).$$

Then, applying  $\mathfrak{D}$  on both sides of the reverse recursion (7.3):

$$\mathbb{R}^{d_i \times m} \ni \mathfrak{D}(V_i) = \sum_{j \in \mathcal{O}_i} [\mathfrak{D}(\nabla_i f_j) \cdot V_j + \nabla_i f_j \cdot \mathfrak{D}(V_j)] = \sum_{j \in \mathcal{O}_i} \left[ \nabla_i f_j \cdot \mathfrak{D}(V_j) + \left( \sum_{l \in \mathcal{I}_j} \nabla_{il}^2 f_j \times \mathfrak{D}(v_l) \right) \cdot V_j \right].$$

Recall that  $v_j = f_j(\mathcal{I}_j)$ , whence from forward-mode differentiation:

$$\mathbb{R}^{1 \times d_j} \ni \mathfrak{D}(v_j) = \sum_{l \in \mathcal{I}_j} \mathfrak{D}(v_l) \cdot \nabla_l f_j, \quad \text{where} \quad \nabla_l f_j = \frac{\partial f_j}{\partial v_l} \in \mathbb{R}^{d_j \times d_l}.$$

Augmenting the assumption (7.2) to include  $\nabla^2 f$ , we conclude that the time complexity of the algorithm below is on par with that of the reverse-mode auto-differentiation in Algorithm 7.5.

---

**Algorithm:** Reverse-mode Hessian-vector product for superposition (Kim et al. 1984, p. 58)

---

**Input:**  $\mathbf{x} \in \mathbb{R}^d$ , direction  $\mathbf{z} \in \mathbb{R}^d$ , basic function class  $\mathcal{F}_0$ , computational graph  $\mathcal{G}$

**Output:** Hessian-vector product  $[\mathfrak{D}(V_1); \dots; \mathfrak{D}(V_d)] \in \mathbb{R}^{d \times m}$

```

1 Run reverse-mode AD in Algorithm 7.5
2 for  $i = 1, \dots, d$  do // forward: initialize directional derivative
3    $\mathfrak{D}(v_i) \leftarrow z_i$ 
4 for  $j = d+1, \dots, d+k+m$  do // forward: accumulate directional derivatives
5    $\mathfrak{D}(v_j) \leftarrow \sum_{l \in \mathcal{I}_j} \mathfrak{D}(v_l) \cdot \nabla_l f_j$  //  $\mathfrak{D}(v_j) \in \mathbb{R}^{1 \times d_j}$ 
6    $\dot{\mathbf{g}}_{:j} \leftarrow \sum_{l \in \mathcal{I}_j} \nabla_{:l}^2 f_j \times \mathfrak{D}(v_l)$  //  $\dot{\mathbf{g}}_{:j} \in \mathbb{R}^{d_i \times d_j}$ 
7 for  $i = 1, \dots, m$  do // backward: initialize Hessian-vector product
8    $\mathfrak{D}(V_{d+k+i}) \leftarrow \mathbf{0}$  //  $\mathfrak{D}(V_{d+k+i}) \in \mathbb{R}^{1 \times m}$ 
9 for  $i = d+k, \dots, 1$  do // backward: accumulate Hessian-vector product
10   $\mathfrak{D}(V_i) \leftarrow \sum_{j \in \mathcal{O}_i} [\nabla_i f_j \cdot \mathfrak{D}(V_j) + \dot{\mathbf{g}}_{:j} \cdot V_j]$  //  $\mathfrak{D}(V_i) \in \mathbb{R}^{d_i \times m}$ 

```

---

We can similarly derive the forward-mode Hessian-vector product:

$$\mathbb{R}^{d \times d_i} \ni \mathfrak{D}(U_i) = \sum_{j \in \mathcal{I}_i} [\mathfrak{D}(U_j) \cdot \nabla_j f_i + U_j \cdot \mathfrak{D}(\nabla_j f_i)] = \sum_{j \in \mathcal{I}_i} \left[ \mathfrak{D}(U_j) \cdot \nabla_j f_i + U_j \cdot \sum_{l \in \mathcal{I}_j} \nabla_{jl}^2 f_i \times \mathfrak{D}(v_l) \right].$$

Running the above algorithms with  $\mathbf{z} \in \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d\}$  recovers the entire Hessian, which is thus  $d$  times more expensive than computing the Hessian-vector product or the gradient.

Needless to say, the same idea extends to computing even higher order of derivatives (Kim et al. 1984, p. 62), including possibly those of nonsmooth functions.

Kim, K. V., Y. E. Nesterov, V. A. Skokov, and B. V. Cherkasskii (1984). “An efficient algorithm for computing derivatives and extremal problems”. *Ekonomika i matematicheskie metody*, vol. 20, no. 2, pp. 309–318.

**Example 7.13: Hessian-vector product in practice**

We mention the following example applications of fast Hessian-vector product.

- Conjugate gradient using Algorithm 7.12 (Kim et al. 1984):

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{s}_t, \quad \text{where}$$

$$\mathbf{s}_t = \nabla f(\mathbf{w}_t) - \frac{\langle \nabla^2 f(\mathbf{w}_t) \mathbf{s}_{t-1}, \nabla f(\mathbf{w}_t) \rangle}{\langle \nabla^2 f(\mathbf{w}_t) \mathbf{s}_{t-1}, \mathbf{s}_{t-1} \rangle} \mathbf{s}_{t-1}$$

$$\eta_t = \operatorname{argmin}_{\eta \geq 0} f(\mathbf{w}_t - \eta \mathbf{s}_t)$$

- Newton’s method using Algorithm 7.12 (Kim et al. 1984), where the update can be reformulated as a quadratic minimization problem, solved through conjugate gradient.
- See also Christianson (1992), Møller (1993), Pearlmutter (1994), and Schraudolph (2002).

Kim, K. V., Y. E. Nesterov, V. A. Skokov, and B. V. Cherkasskii (1984). “An efficient algorithm for computing derivatives and extremal problems”. *Ekonomika i matematicheskie metody*, vol. 20, no. 2, pp. 309–318.

Christianson, B. (1992). “Automatic Hessians by reverse accumulation”. *IMA Journal of Numerical Analysis*, vol. 12, no. 2, pp. 135–150.

Møller, M. (1993). “Exact Calculation of the Product of the Hessian Matrix of Feed-Forward Network Error Functions and a Vector in  $O(N)$  Time”. Tech. rep. DAIMI Report Series, 22(432).

Pearlmutter, B. A. (1994). “Fast Exact Multiplication by the Hessian”. *Neural Computation*, vol. 6, no. 1, pp. 147–160.

Schraudolph, N. N. (2002). “Fast Curvature Matrix-Vector Products for Second-Order Gradient Descent”. *Neural Computation*, vol. 14, no. 7, pp. 1723–1738.

#### Remark 7.14: Approximating the diagonal of Hessian (e.g., LeCun et al. 1989)

Differentiating again the recursion (7.3) we obtain

$$\mathbb{R}^{d_i \times d_i \times m} \ni \frac{\partial^2 \mathbf{y}}{\partial v_i^2} = \sum_{j \in \mathcal{O}_i} \left[ \nabla_i^2 f_j \frac{\partial \mathbf{y}}{\partial v_j} + \nabla_i f_j \frac{\partial^2 \mathbf{y}}{\partial v_j \partial v_i} \right] \approx \sum_{j \in \mathcal{O}_i} \nabla_i^2 f_j \frac{\partial \mathbf{y}}{\partial v_j}. \quad (7.5)$$

Clearly, the right-hand side is only an approximation of the diagonal of the Hessian and can be computed at the same cost as the gradient. LeCun et al. (1989) used the Taylor expansion

$$\Delta f(\mathbf{w}) \approx \langle \nabla f(\mathbf{w}), \Delta \mathbf{w} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{w}) \Delta \mathbf{w}, \Delta \mathbf{w} \rangle$$

to assess the saliency of components of  $\mathbf{w}$ . At a local optimum,  $\nabla f(\mathbf{w}) \approx \mathbf{0}$ , so we only need to approximate the quadratic term. By dropping the off-diagonal terms as in (7.5),

LeCun, Y., J. Denker, and S. Solla (1989). “Optimal Brain Damage”. In: *Advances in Neural Information Processing Systems 2*.