

# VECTOR DIFFERENTIATION RULES

Derivatives with respect to a vector come up often in a multitude of areas, such as constrained optimisation, adaptive filtering, and machine learning. We begin by defining the simple case of a scalar differentiated by a vector.

**DEFINITION 0.1 (gradient).** Let  $\mathbf{x} = (x_1, \dots, x_n)$  be a column vector and let  $f(\mathbf{x}): \mathbb{R}^n \rightarrow \mathbb{R}$  be a function that maps to a scalar. The derivative of  $f$  w.r.t  $\mathbf{x}$ , also known as *gradient*, is defined as:

$$\nabla_{\mathbf{x}} f := \frac{\partial f}{\partial \mathbf{x}} := \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \dots \quad \frac{\partial f}{\partial x_n} \right]^\top \quad (0.1)$$

Both notations in Eq. (0.1) are acceptable for the gradient

**EXAMPLE 0.1.** Find the gradient of the function  $f(x_1, x_2, x_3) = x_1 + 3x_2 + 2x_3$ .

**SOLUTION 0.1.**

$$\nabla_{\mathbf{x}} f = \left[ \frac{\partial(x_1+3x_2+2x_3)}{\partial x_1} \quad \frac{\partial(x_1+3x_2+2x_3)}{\partial x_2} \quad \frac{\partial(x_1+3x_2+2x_3)}{\partial x_3} \right]^\top = [1 \quad 3 \quad 2]^\top$$

□

We already notice a property of vector differentiation in the example above; if  $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$ , where  $\mathbf{a}^\top = [1 \quad 3 \quad 2]$  is a coefficient row vector and  $\mathbf{x} = [x_1 \quad x_2 \quad x_3]^\top$  a column vector, then  $\nabla_{\mathbf{x}} f = \mathbf{a}$ .

For reference, the derivative of vector w.r.t. a vector is also defined just to highlight its difference with the derivative of a scalar w.r.t. a vector.

**DEFINITION 0.2 (Jacobian matrix).** Let  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{f}: \mathbb{R}^m \rightarrow \mathbb{R}^n$  be a function that returns a  $n \times 1$  vector, i.e.  $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}) \quad \dots \quad f_n(\mathbf{x})]^\top$ . Then the derivative of vector  $\mathbf{f}(\mathbf{x})$  w.r.t.  $\mathbf{x}$  is called *Jacobian matrix* and is defined as [1]

$$\mathbf{J}(\mathbf{x}) = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_m} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_2(\mathbf{x})}{\partial x_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n(\mathbf{x})}{\partial x_1} & \frac{\partial f_n(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_n(\mathbf{x})}{\partial x_m} \end{bmatrix} \quad (0.2)$$

The determinant of the Jacobian matrix is called Jacobian determinant or *Jacobian* for short. Therefore each row  $k$  contains the derivative of the scalar function  $f_k(\cdot)$  with respect to the elements in  $\mathbf{x}$ .

**EXAMPLE 0.2.** Compute the Jacobian matrix of the transformation  $T(u, v) = [u \quad v \quad u^v]^\top$ ,  $u > 0$  [2].

**SOLUTION 0.2.** Using the notation from the definition we compute each row at a time.

$$\begin{aligned} f_1(u, v) = u &\Rightarrow \frac{\partial f_1(u, v)}{\partial u} = 1, \quad \frac{\partial f_1(u, v)}{\partial v} = 0 \\ f_2(u, v) = v &\Rightarrow \frac{\partial f_2(u, v)}{\partial u} = 0, \quad \frac{\partial f_2(u, v)}{\partial v} = 1 \\ f_3(u, v) = u^v &\Rightarrow \frac{\partial f_3(u, v)}{\partial u} = v u^{v-1}, \quad \frac{\partial f_3(u, v)}{\partial v} = u^v \ln u \\ \therefore \mathbf{J}(u, v) &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ v u^{v-1} & u^v \ln u \end{bmatrix} \end{aligned}$$

□

Now we can derive and provide the derivatives of some common scalar ((i), (iii), (iv)) and one vector ((ii)) expressions w.r.t. a vector.

**LEMMA 0.1 (vector differentiation basic properties).** Let  $\mathbf{x}, \mathbf{a} \in \mathbb{R}^n$  be two column vectors where  $\mathbf{a}$  is not a

function of  $\mathbf{x}$  and  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a real matrix. Then:

$$(i) \quad \frac{\partial(\mathbf{a}^\top \mathbf{x})}{\partial \mathbf{x}} = \frac{\partial(\mathbf{x}^\top \mathbf{a})}{\partial \mathbf{x}} = \mathbf{a} \quad (0.3)$$

$$(ii) \quad \frac{\partial(\mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A} \quad (0.4)$$

$$(iii) \quad \frac{\partial(\mathbf{x}^\top \mathbf{A}^\top)}{\partial \mathbf{x}} = \mathbf{A}^\top, \quad \text{if } m = n \quad (0.5)$$

$$(iv) \quad \frac{\partial(\mathbf{x}^\top \mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}, \quad \text{if } m = n \quad (0.6)$$

*Proof.*

(i) From the dot product's definition:

$$\frac{\partial(\mathbf{a}^\top \mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\sum_{i=1}^n a_i x_i}{\partial x_1} \\ \frac{\sum_{i=1}^n a_i x_i}{\partial x_2} \\ \vdots \\ \frac{\sum_{i=1}^n a_i x_i}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \mathbf{a}$$

(ii) If we denote  $\mathbf{a}_1^\top, \dots, \mathbf{a}_m^\top$  the rows of  $\mathbf{A}$  expressed as column vectors, then product  $\mathbf{A}\mathbf{x}$  is written as:

$$\mathbf{A}\mathbf{x} = \begin{bmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \mathbf{a}_2^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_m^\top \mathbf{x} \end{bmatrix}$$

We apply definition Eq. (0.1) on each row, since each row is a scalar:

$$\therefore \frac{\partial(\mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial(\mathbf{a}_1^\top \mathbf{x})}{\partial \mathbf{x}} \stackrel{(0.1)}{=} \mathbf{a}_1 \\ \frac{\partial(\mathbf{a}_2^\top \mathbf{x})}{\partial \mathbf{x}} \stackrel{(0.1)}{=} \mathbf{a}_2 \\ \vdots \\ \frac{\partial(\mathbf{a}_m^\top \mathbf{x})}{\partial \mathbf{x}} \stackrel{(0.1)}{=} \mathbf{a}_m \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_m \end{bmatrix} = \mathbf{A}$$

(iii) Left as exercise.

(iv) The  $i$ -th element of the product  $\mathbf{A}\mathbf{x}$ , which is a vector, is written with the index notation as follows.

$$(\mathbf{A}\mathbf{x})_i = \sum_{j=1}^n A_{ij} x_j$$

The dot product of  $\mathbf{x}$ ,  $\mathbf{A}\mathbf{x}$  is written as:

$$\mathbf{x}^\top \mathbf{A}\mathbf{x} = \sum_i^n x_i \sum_{j=1}^n A_{ij} x_j$$

Applying the definition of the gradient (Eq. (0.1)) to the dot product:

$$\begin{aligned} \frac{\partial(\mathbf{x}^\top \mathbf{A}\mathbf{x})}{\partial \mathbf{x}} &= \frac{\partial(\sum_i^n x_i \sum_{j=1}^n A_{ij} x_j)}{\partial \mathbf{x}} \\ &= \left[ \frac{\partial(\sum_{i=1}^n x_i \sum_{j=1}^n A_{ij} x_j)}{\partial x_1} \quad \dots \quad \frac{\partial(\sum_{i=1}^n x_i \sum_{j=1}^n A_{ij} x_j)}{\partial x_n} \right]^\top \end{aligned}$$

Using the product rule on the first element:

$$\begin{aligned}
\frac{\sum_{i=1}^n x_i \sum_{j=1}^n A_{ij} x_j}{\partial x_1} &= \cancel{\frac{\partial \sum_{i=1}^n x_i}{\partial x_1}} \overset{1, i=1}{\text{else } 0} \sum_{j=1}^n A_{ij} x_j + \sum_{i=1}^n x_i \cancel{\frac{\partial \sum_{j=1}^n A_{ij} x_j}{\partial x_1}} \overset{A_{ij}, j=1}{\text{else } 0} \\
&= \sum_{j=1}^n A_{1j} x_j + \sum_{i=1}^n x_i A_{i1} \\
&= \mathbf{a}_{1:} \mathbf{x} + \mathbf{a}_{:1} \mathbf{x} = (\mathbf{a}_{1:} + \mathbf{a}_{:1}) \mathbf{x}
\end{aligned}$$

, where  $\mathbf{a}_{1:}$  denotes the first row of  $\mathbf{A}$  and  $\mathbf{a}_{:1}$  its first column (Matlab notation). Applying the result to the remaining indexes  $2, \dots, n$ , we can rewrite the gradient as:

$$\frac{\partial(\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = [\mathbf{a}_{1:} + \mathbf{a}_{:1} \quad \dots \quad \mathbf{a}_{n:} + \mathbf{a}_{:n}] \mathbf{x} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

□

Some other properties that can be readily derived from the basic ones are listed below.

**LEMMA 0.2 (vector differentiation follow-up properties).** If  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is square and symmetric ( $\mathbf{A} = \mathbf{A}^\top$ ) and  $\mathbf{x} \in \mathbb{R}^n$ , then

$$\frac{\partial(\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{A} \mathbf{x} \tag{0.7}$$

For  $\mathbf{A} = \mathbf{I}$ , we can derive the vector derivative of the squared norm-2:

$$\frac{\partial(\mathbf{x}^\top \mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \|\mathbf{x}\|^2}{\partial \mathbf{x}} = 2\mathbf{x} \tag{0.8}$$

Finally, we define the chain rule for vector functions as it's expressed in a slightly different order than the chain rule of scalars.

**LEMMA 0.3 (chain rule for vector function).** Let  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^r$ ,  $\mathbf{z} \in \mathbb{R}^m$ , where  $\mathbf{z} = \mathbf{z}(\mathbf{y})$  and  $\mathbf{y} = \mathbf{y}(\mathbf{x})$ . Then

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \tag{0.9}$$

*Proof.*

Proof is found in [3].

□

## References

- [1] H. B. Nielsen, *Introduction to vector and matrix differentiation*, 2012. [Online]. Available: [https://absalon.instructure.com/files/1853451/download?download\\_frd=1](https://absalon.instructure.com/files/1853451/download?download_frd=1).
- [2] J. Ruan, *Jacobians and their applications*. [Online]. Available: <https://www.projectrhea.org/rhea/index.php/Jacobian>.
- [3] W. Zhu, *Matrix & vector basic linear algebra & calculus*. [Online]. Available: <http://www.ams.sunysb.edu/~zhu/ams571/matrixvector.pdf>.