# Vector Calculus

## Gradient

For a scalar valued function $f$ with vector valued input $\mathbf{x} \in \mathbb{R}^n$ , the derivative is this row vector of shape $1 \times n$ .

$$\frac{df}{d\mathbf{x}} = \nabla_{\mathbf{x}}f = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} & \dfrac{\partial f}{\partial x_2} & \cdots & \dfrac{\partial f}{\partial x_n} \end{bmatrix}$$

This is the vector along with rate of increase of $f$ (wrt norm of displacement) is highest.

## Jacobian for vectors

For a vector valued function $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$ with vector valued input $\mathbf{x} \in \mathbb{R}^n$ , the derivative is this matrix of shape $m \times n$ called the Jacobian matrix, or simply the Jacobian :

$$\frac{d\mathbf{f}}{d\mathbf{x}} = \nabla_{\mathbf{x}}\mathbf{f} = \mathbf{J} \quad \text{where} \quad J_{ij} = \frac{\partial f_i}{\partial x_j}$$

## Jacobian for Matrices

For a matrix valued function $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^{p \times q}$ with matrix valued $\mathbf{x} \in \mathbb{R}^{m \times n}$ , the derivative is again the Jacobian $\mathbf{J}$ which is now a tensor of shape $p \times q \times m \times n$ , which can't be written on paper. $\mathbf{J}$ is given by its elements instead as :

$$J_{ijkl} = \frac{\partial f_{ij}}{\partial x_{kl}}$$

But this approach is impractical and what is usually done is that the matrices $\mathbf{f}, \mathbf{x}$ are flattened out into vectors of size $pq$ and $mn$ , and then the Jacobian matrix with respect to these vectors is calculated.

## Useful identities for calculating Jacobians

$$\frac{\partial}{\partial X} f(X)^\top = \left(\frac{\partial f(X)}{\partial X}\right)^\top$$

$$\frac{\partial}{\partial X} \text{tr}(f(X)) = \text{tr}\left(\frac{\partial f(X)}{\partial X}\right)$$

$$\frac{\partial}{\partial X} \det(f(X)) = \det(f(X))\text{tr}\left(f(X)^{-1}\frac{\partial f(X)}{\partial X}\right)$$

$$\frac{\partial}{\partial X} f(X)^{-1} = -f(X)^{-1}\frac{\partial f(X)}{\partial X}f(X)^{-1}$$

$$\frac{\partial a^\top X^{-1} b}{\partial X} = -(X^{-1})^\top a b^\top (X^{-1})^\top$$

$$\frac{\partial x^\top a}{\partial x} = a^\top$$

$$\frac{\partial a^\top x}{\partial x} = a^\top$$

$$\frac{\partial a^\top X b}{\partial X} = a b^\top$$

$$\frac{\partial x^\top B x}{\partial x} = x^\top (B + B^\top)$$

$$\frac{\partial}{\partial s}(x - As)^{\top} W(x - As) = -2(x - As)^{\top} W A \quad \text{for symmetric } W$$

## Automatic differentiation

Sometimes the derivative of a complex function (a function composed of many small functions) is even more complex than the function itself. Suppose $f$ is such a function composed of $n$ simple functions, as $f = f_n \circ f_{n-1} \ldots f_2 \circ f_1$ and $f'$ is composed of $m > n$ simple functions. Then instead of computing the derivative directly, we can keep track of the values of $k$ succesive functions applied on the current input as $x_k = f_k \circ f_{k-1} \circ f_{k-2} \ldots f_2 \circ f_1(x)$ . Then using $\frac{dx_k}{dx} = f'_k(x_{k-1})\frac{dx_{k-1}}{dx}$ , we calculate the derivatives of all $x_k$ and update them using Euler's method (a numerical method for simulating differential equations), then in the last iteration of every update, we arrive at $f'(x) = \frac{dx_n}{dx}$ . Here, we've done only $n$ simple computations rather than $m$ computations done when using the explicit form of $f'$ . Note that it's important that you update every $x_i$ only after you have used it to calculate $f'(x_{k-1})$ . Also in practice, the equation you would use is not $\frac{dx_k}{dx} = f'_k(x_{k-1})\frac{dx_{k-1}}{dx}$ , but just $\Delta x_k = f'(x_{k-1})\Delta x_{k-1}$ .

This can also be generalised to functions composed of vector valued simple functions with vector valued inputs.
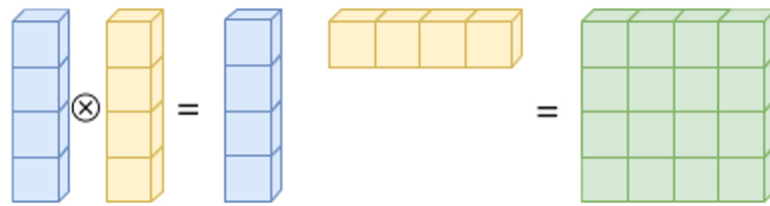
## Second order derivative wrt  a vector

For a scalar valued function $f(\mathbf{x})$ with inputs $\mathbf{x} \in \mathbb{R}^n$ , the first order derivative is $\frac{df}{d\mathbf{x}} = \big[\frac{\partial f}{\partial x_i}\big]_i^T$ where the square bracket means "stack one below other for different $i$" , i.e. create a column vector with such and such entries. Notice that the structure in which the vectors are represented doesn't matter all that much, the containment does. (if you have ever used the numpy library in python, you
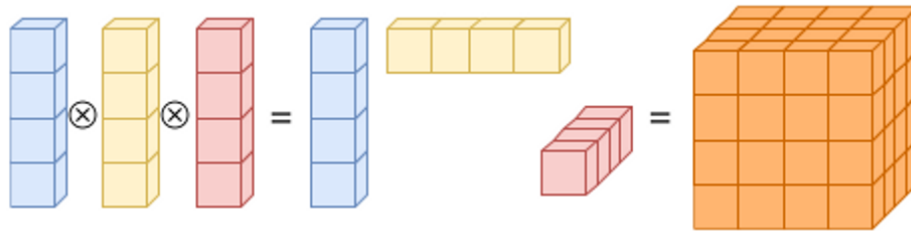
know what I'm talking about) So now, the 2nd derivative is $\frac{d}{d\mathbf{x}}\left(\left[\frac{\partial f}{\partial x_i}\right]_i\right) = \left[\frac{\partial^2 f}{\partial x_i \partial x_j}\right]_{ij} = \mathbf{H}$ , also known as the Hessian.

You can also think of this as $\nabla^T \nabla f$ where $\nabla = \left[\begin{array}{cccc}\frac{\partial}{\partial x_1} & \frac{\partial}{\partial x_2} & \cdots & \frac{\partial}{\partial x_n}\end{array}\right]$ .

To go to the third derivative, we need a generalised notion of matrices that doesn't depend on the orientation and other trivialities. thus we use tensors and the outer product, which behaves like matrix multiplication :



(a) Given a vector $\boldsymbol{\delta} \in \mathbf{R}^4$, we obtain the outer product $\boldsymbol{\delta}^2 := \boldsymbol{\delta} \otimes \boldsymbol{\delta} = \boldsymbol{\delta}\boldsymbol{\delta}^\top \in \mathbf{R}^{4\times 4}$ as a matrix.



(b) An outer product $\boldsymbol{\delta}^3 := \boldsymbol{\delta} \otimes \boldsymbol{\delta} \otimes \boldsymbol{\delta} \in \mathbf{R}^{4\times 4 \times 4}$ results in a third-order tensor ("three-dimensional matrix"), i.e., an array with three indexes.

Now, we define $\nabla_{\mathbf{x}}^k$ to be the $k$-th total derivative wrt $\mathbf{x}$ , which is given by taking $k$ outer products of the nabla operator. Basically, $\nabla_{\mathbf{x}}^k = \nabla \otimes \nabla \otimes \nabla \otimes \ldots \nabla$ (k times).

## Multivariate Taylor series

For an analytical function $f(\mathbf{x})$ , if $\mathbf{x} = \mathbf{x_0} + \delta$ , then we have the taylor series :

$$f(\mathbf{x}) = \sum_{k=0}^{\infty} \frac{\nabla^k f(\mathbf{x_0})}{k!} \cdot \delta^{\mathbf{k}}$$

Here $\delta^{\mathbf{k}} = \delta \otimes \delta \otimes \delta \otimes \ldots \delta$ is the outer product of $\delta$ taken $k$ times. we are taking the inner product of $\nabla^k f(\mathbf{x_0})$ and $\delta^k$ as

$$\nabla^k f(\mathbf{x_0}) \cdot \delta^k = \sum_{i_1, i_2 \ldots i_n} (\nabla^k f(\mathbf{x_0}))[i_1, i_2 \ldots i_n]\, \delta^k[i_1, i_2, i_3 \ldots i_n] =$$

$$\sum_{i_1, i_2 \ldots i_n} (\nabla^k f(\mathbf{x_0}))[i_1, i_2 \ldots i_n]\, \delta_{i_1}\, \delta_{i_2}\, \delta_{i_3} \ldots \delta_{i_n} =$$

$$\sum_{i_1, i_2 \ldots i_n} (\frac{\partial}{\partial x_{i_1}} \frac{\partial}{\partial x_{i_2}} \frac{\partial}{\partial x_{i_3}} \cdots \frac{\partial}{\partial x_{i_n}} f)_{(\mathbf{x_0})}\, \delta_{i_1}\, \delta_{i_2}\, \delta_{i_3} \ldots \delta_{i_n}$$