

Differentiating Vector- and Matrix-Valued Functions

Szymon Rusinkiewicz
COS 302, Fall 2020



Generalizing Functions...

Functions of scalars, vectors, matrices ... *returning* scalars, vectors, matrices

- Function of a scalar, returning a scalar: $\mathbb{R} \rightarrow \mathbb{R}$
 - Example: $f(x) = ax + b$
- Function of a scalar, returning a vector: $\mathbb{R} \rightarrow \mathbb{R}^n$
 - Example: $f(x) = x\mathbf{v}$
- Function of a vector, returning a scalar: $\mathbb{R}^n \rightarrow \mathbb{R}$
 - Example: $f(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$
- Function of a vector, returning a vector: $\mathbb{R}^n \rightarrow \mathbb{R}^n$
 - Example: $f(\mathbf{x}) = \mathbf{M}\mathbf{x}$
- Function of a vector, returning a matrix: $\mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$
 - Example: $F(\mathbf{x}) = \mathbf{x}\mathbf{x}^\top$
- Many other possibilities: function of a matrix, etc.

Generalizing Functions... and Taking Their Derivatives

- Function of a scalar, returning a scalar: $\mathbb{R} \rightarrow \mathbb{R}$
 - Example: $f(x) = ax + b$ \rightarrow Ordinary derivative $\frac{df}{dx} : \mathbb{R} \rightarrow \mathbb{R}$
- Function of a scalar, returning a vector: $\mathbb{R} \rightarrow \mathbb{R}^n$
 - Example: $f(x) = x\mathbf{v}$ \rightarrow Vector-valued derivative $\frac{df}{dx} : \mathbb{R} \rightarrow \mathbb{R}^n$
- Function of a vector, returning a scalar: $\mathbb{R}^n \rightarrow \mathbb{R}$
 - Example: $f(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$ \rightarrow Gradient $\nabla f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^{1 \times n}$
- Function of a vector, returning a vector: $\mathbb{R}^n \rightarrow \mathbb{R}^n$
 - Example: $f(\mathbf{x}) = M\mathbf{x}$ \rightarrow Jacobian $\nabla f(\mathbf{x})$ or $J(f(\mathbf{x})) : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$
- Function of a vector, returning a matrix: $\mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$
 - Example: $F(\mathbf{x}) = \mathbf{x}\mathbf{x}^\top$ \rightarrow Generalized Jacobian $\nabla F(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n \times n}$

Generalizing Functions... and Taking Their Derivatives

In general, if f is a function

$$f : \mathbb{R}^{(\text{input shape})} \rightarrow \mathbb{R}^{(\text{output shape})}$$

then its generalized derivative will be a function

$$\nabla f : \mathbb{R}^{(\text{input shape})} \rightarrow \mathbb{R}^{(\text{output shape}) \times (\text{input shape})}$$

where the extra dimensions on output correspond to taking partial derivatives with respect to all the input dimensions.

Tensors

So what if we end up with e.g. an $n \times n \times n$ object?

- *Tensors* are multidimensional generalizations of scalars, vectors, matrices.
- For our purposes, represented as multidimensional arrays of numbers.
- The number of indices in the shape can be called the *order*, *degree*, (confusingly) *dimension*, or (even more confusingly) *rank* of the tensor.
 - Example: Take a function of a vector, returning a matrix, and differentiate it. The resulting $n \times n \times n$ beastie is a degree-3 or 3rd-order tensor.

Tensors in Python

- NumPy arrays can represent tensors
 - `A = np.zeros((5, 6, 7))` → `A.shape == (5, 6, 7)`
- Transpose can take a permutation of dimensions
 - `B = np.transpose(A, (2, 0, 1))` → `B.shape == (7, 5, 6)`
- Careful if using `np.matmul` or `np.dot` for tensor multiplication — `np.tensordot` lets you explicitly specify axes to sum over
 - `C = np.tensordot(A, B, (2, 0))` → `C.shape == (5, 6, 5, 6)`
 - `D = np.tensordot(A, B, ([2,1], [0,2]))` → `D.shape == (5, 5)`

Tensors and Differentiation

$$\nabla f : \mathbb{R}^{(\text{input shape})} \rightarrow \mathbb{R}^{(\text{output shape}) \times (\text{input shape})}$$

- If you take more advanced math, you'll learn that the “dimensions” of tensors behave in two different ways: *covariant* and *contravariant*.
- We won't go into that here, except to note that the dimensions arising from differentiation always behave “transposed”.
 - For example, the gradient of a scalar function of a vector is a *row* vector.
 - Intimate connection to *directional derivatives*: multiplying a gradient by a *direction* \mathbf{d} (an object of the input shape) gives you derivative of the output in that direction:

$$D_{\mathbf{d}} f(\mathbf{x}) \text{ or } \nabla_{\mathbf{d}} f(\mathbf{x}) = \nabla f(\mathbf{x}) \mathbf{d}$$

where the last dimension(s) of ∇f are dotted against \mathbf{d} .

Preliminaries

Before we get into specific examples of these generalized derivatives, let's review which rules from single-variable calculus still work:

- Derivative of a *constant* of any shape is 0
- Derivative of the variable with respect to which we're differentiating is 1, or the *identity* of the appropriate shape
- Derivative of a *sum* is the sum of derivatives
- Derivative of a *scalar multiple* is the constant times the derivative
- *Chain rule* works, but order might matter: $\frac{d}{dx} f(g(x)) = f'(g(x)) g'(x)$
- *Product rule* requires care about dimensions and transposes (stay tuned!)

Function of a Scalar, Returning a Vector

Simple...

$$f(x) = \begin{bmatrix} 3x + 42 \\ \sin x \end{bmatrix}$$
$$\frac{df}{dx} = \begin{bmatrix} 3 \\ \cos x \end{bmatrix}$$

Can also consider functions written as scalar/vector products:

$$f(x) = x\mathbf{v}$$
$$\frac{df}{dx} = \begin{bmatrix} \frac{d}{dx}(x v_1) \\ \frac{d}{dx}(x v_2) \\ \vdots \end{bmatrix} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \end{bmatrix} = \mathbf{v}$$

Function of a Vector, Returning a Scalar

This is the ordinary *gradient*, which is a *row* vector of partial derivatives:

$$f\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = x_1^3 + x_1x_2 + 42x_3$$
$$\begin{aligned}\nabla f &= \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \frac{\partial f}{\partial x_3} \right] \\ &= [3x_1^2 + x_2 \quad x_1 \quad 42]\end{aligned}$$

Directional Derivative

$$\nabla f = [3x_1^2 + x_2 \quad x_1 \quad 42]$$

How does f change with an infinitesimal step in direction $\mathbf{d} = \begin{bmatrix} 0.6 \\ 0 \\ 0.8 \end{bmatrix}$?

$$\nabla_{\mathbf{d}} f = [3x_1^2 + x_2 \quad x_1 \quad 42] \begin{bmatrix} 0.6 \\ 0 \\ 0.8 \end{bmatrix} = 1.8x_1^2 + 0.6x_2 + 33.6$$

This is a *scalar*—the same shape as the output of f .

Directional Derivative

- What if we had done this in the previous case, where we had a function of a *scalar*, returning a *vector*?

$$f(x) = x\mathbf{v}$$
$$\frac{df}{dx} = \mathbf{v}$$

- Multiplying by a (scalar) infinitesimal step in x in “direction” 1, we get just \mathbf{v} .
- This is a (column) vector — the same shape as the output of f .

Function of a Vector, Returning a Scalar

Let's try a dot product!

$$\begin{aligned}f(\mathbf{x}) &= \mathbf{v} \cdot \mathbf{x} = \sum_i v_i x_i \\ \nabla f &= \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \right] \\ &= [v_1 \quad v_2 \quad \cdots] \\ &= \mathbf{v}^\top\end{aligned}$$

But $\mathbf{v} \cdot \mathbf{x} = \mathbf{v}^\top \mathbf{x} = \mathbf{x}^\top \mathbf{v}$, so we have the following:

$$\nabla(\mathbf{v}^\top \mathbf{x}) = \mathbf{v}^\top \quad \text{and} \quad \nabla(\mathbf{x}^\top \mathbf{v}) = \mathbf{v}^\top$$

Function of a Vector, Returning a Scalar

Next interesting case:

$$f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{x} = \mathbf{x}^\top \mathbf{x} = \|\mathbf{x}\|^2 = \sum_i x_i^2$$

$$\begin{aligned}\nabla f &= \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \cdots \end{bmatrix} \\ &= \begin{bmatrix} 2x_1 & 2x_2 & \cdots \end{bmatrix} \\ &= 2\mathbf{x}^\top\end{aligned}$$

Note the analogy to $\frac{d}{dx}x^2 = 2x$, but we need the transpose to get the output shape right.

Function of a Vector, Returning a Scalar

Even more interesting:

$$f(\mathbf{x}) = \|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}} = \sqrt{\sum x_i^2}$$

Applying the chain rule:

$$\begin{aligned}\nabla f &= \frac{1}{2}(\mathbf{x}^\top \mathbf{x})^{-\frac{1}{2}} \nabla(\mathbf{x}^\top \mathbf{x}) \\ &= \frac{2\mathbf{x}^\top}{2\sqrt{\mathbf{x}^\top \mathbf{x}}} \\ &= \frac{\mathbf{x}^\top}{\|\mathbf{x}\|}\end{aligned}$$

Directional Derivative

$$\nabla \|\mathbf{x}\| = \frac{\mathbf{x}^\top}{\|\mathbf{x}\|}$$

As \mathbf{x} changes by an infinitesimal step in direction \mathbf{d} ,

$$\nabla_{\mathbf{d}} \|\mathbf{x}\| = \frac{\mathbf{x}}{\|\mathbf{x}\|} \cdot \mathbf{d}$$

Intuitive: if \mathbf{d} is in the direction of \mathbf{x} , change in $\|\mathbf{x}\|$ is 1 times the step size, etc.

Function of a Vector, Returning a Vector

Let's move on to a Jacobian:

$$\begin{aligned}f(\mathbf{x}) &= \mathbf{M}\mathbf{x} \\ \nabla f &= \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \right]\end{aligned}$$

The only terms in $\mathbf{M}\mathbf{x}$ involving x_i come from the i^{th} column of \mathbf{M} , so:

$$\frac{\partial f}{\partial x_1} = \begin{bmatrix} M_{11} \\ M_{21} \\ \vdots \end{bmatrix}, \quad \frac{\partial f}{\partial x_2} = \begin{bmatrix} M_{12} \\ M_{22} \\ \vdots \end{bmatrix}, \text{ etc.}$$

Function of a Vector, Returning a Vector

- Stitching everything together,

$$\begin{aligned}\nabla(\mathbf{M}\mathbf{x}) &= \left[\begin{array}{c} \left(\begin{array}{c} M_{11} \\ M_{21} \\ \vdots \end{array} \right) & \left(\begin{array}{c} M_{12} \\ M_{22} \\ \vdots \end{array} \right) & \dots \end{array} \right] \\ &= \mathbf{M}\end{aligned}$$

- Special case: $\nabla(\mathbf{x}) = \nabla(\mathbf{I}\mathbf{x}) = \mathbf{I}$
- This reinforces our intuition that differentiating any constant thing times \mathbf{x} gives just that constant, whether it's a scalar, vector, matrix, etc.

Function of a Vector, Returning a Matrix

$$F(\mathbf{x}) = \mathbf{x}\mathbf{v}^T = \begin{bmatrix} x_1v_1 & x_1v_2 & \cdots \\ x_2v_1 & x_2v_2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Just apply the rules, and watch the tensor appear!

$$\begin{aligned} \nabla F(\mathbf{x}) &= \begin{bmatrix} \frac{\partial F}{\partial x_1} & \frac{\partial F}{\partial x_2} & \cdots \end{bmatrix} \\ &= \begin{bmatrix} \begin{bmatrix} v_1 & v_2 & \cdots \\ 0 & 0 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} & \begin{bmatrix} 0 & 0 & \cdots \\ v_1 & v_2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} & \cdots \end{bmatrix} \end{aligned}$$

Function of a Vector, Returning a Matrix

$$\nabla F(\mathbf{x}) = \left[\begin{array}{c} \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \mathbf{v}^\top \\ \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \mathbf{v}^\top \\ \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix} \mathbf{v}^\top \\ \dots \end{array} \right]$$

At this point, it might be tempting to factor out the \mathbf{v}^\top . But **be careful!**

This object is an $n \times n \times n$ tensor, and if you factor it into an $n \times 1 \times n$ tensor and a $1 \times n$ vector, you have to remember which dimensions should be multiplied!

Generalizing Product Rule

- As we just saw, tensor multiplication can get confusing. This complicates cleanly stating a generalized product rule.
- But, let's derive a rule for vector-vector (dot) products:

$$\nabla(\mathbf{v} \cdot \mathbf{w}) = \nabla \left(\sum_i v_i w_i \right)$$

where \mathbf{v} and \mathbf{w} are both potentially functions of \mathbf{x} .

- Writing out the partial derivatives,

$$\nabla(\mathbf{v} \cdot \mathbf{w}) = \left[\frac{\partial(\sum v_i w_i)}{\partial x_1} \quad \frac{\partial(\sum v_i w_i)}{\partial x_2} \quad \dots \right]$$

Generalizing Product Rule

- Because v_i and w_i are just *scalars*, the product rule works normally:

$$\frac{\partial (\sum v_i w_i)}{\partial x_1} = \sum_i \left(v_i \frac{\partial w_i}{\partial x_1} + \frac{\partial v_i}{\partial x_1} w_i \right)$$

- Applying the distributive rule, we get

$$\nabla(\mathbf{v} \cdot \mathbf{w}) = \mathbf{v} \cdot (\nabla \mathbf{w}) + (\nabla \mathbf{v}) \cdot \mathbf{w}$$

- Or, in matrix notation,

$$\nabla(\mathbf{v}^T \mathbf{w}) = \mathbf{v}^T \nabla \mathbf{w} + \mathbf{w}^T \nabla \mathbf{v}$$

Bilinear Form

Let's apply our newly-derived knowledge!

$$\begin{aligned}f(\mathbf{x}) &= \mathbf{x}^\top \mathbf{M} \mathbf{x} \\ \nabla f &= \mathbf{x}^\top \nabla(\mathbf{M} \mathbf{x}) + (\mathbf{M} \mathbf{x})^\top \nabla \mathbf{x} \\ &= \mathbf{x}^\top \mathbf{M} + \mathbf{x}^\top \mathbf{M}^\top \mathbf{I} \\ &= \mathbf{x}^\top (\mathbf{M} + \mathbf{M}^\top)\end{aligned}$$

Note the similarity to $\frac{d}{dx}(ax^2) = 2ax$.

The Grand Finale: Least Squares

- We've mentioned before that our methods for solving overdetermined linear systems of the form $\mathbf{Ax} = \mathbf{b}$ minimize a least-squares residual:

$$\arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|^2$$

- Let's apply the methods we've learned to find the \mathbf{x} that minimizes this, by taking the derivative (gradient) and setting it equal to 0.

The Grand Finale: Least Squares

- Applying the chain rule:

$$\nabla \|Ax - b\|^2 = 2\|Ax - b\| \nabla \|Ax - b\|$$

- And again (order matters!):

$$= 2\|Ax - b\| \frac{(Ax - b)^\top}{\|Ax - b\|} \nabla (Ax - b)$$

- And computing a final Jacobian:

$$= 2(Ax - b)^\top A$$

The Grand Finale: Least Squares

- To check, let's derive this a different way:

$$\begin{aligned}\|Ax - b\|^2 &= (Ax - b)^\top (Ax - b) \\ &= x^\top A^\top Ax - x^\top A^\top b - b^\top Ax + b^\top b \\ \nabla \|Ax - b\|^2 &= x^\top (A^\top A + (A^\top A)^\top) - b^\top A - b^\top A + 0 \\ &= 2x^\top A^\top A - 2b^\top A \\ &= 2(x^\top A^\top - b^\top)A \\ &= 2(Ax - b)^\top A\end{aligned}$$

The Grand Finale: Least Squares

- Setting the gradient equal to a row vector of zeros:

$$2(\mathbf{Ax} - \mathbf{b})^\top \mathbf{A} = \mathbf{0}^\top$$

- Transposing and dividing by 2:

$$\mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}) = \mathbf{0}$$

- And finally, rearranging:

$$\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b}$$