# Lecture 4 Part 1: Gradients and Inner Products in Other Vector Spaces

## MIT 18.S096 Matrix Calculus For Machine Learning And Beyond

*March 5, 2024*

## Table of contents

## Riesz representation theorem

A *Hilbert space* is a continuous vector space with an inner (dot/scalar) product defined. For $\mathbb{R}^n$ (i.e., column vectors), we usually define the inner product as $\boldsymbol{x} \cdot \boldsymbol{y} = \boldsymbol{x}^T \boldsymbol{y}$. For $\mathbb{R}^{n \times m}$ (i.e., matrices), we usually define the inner product as $\text{sum}(\boldsymbol{A} \odot \boldsymbol{B}) = (\text{vec}(A))^T (\text{vec } B) = \text{tr}(A^T B)$. Three properties of a valid inner product:

1. Symmetric: $x \cdot y = y \cdot x$

2. Linear: $x \cdot (\alpha y + \beta z) = \alpha(x \cdot y) + \beta(x \cdot z)$

3. Non-negative: $x \cdot x = \|x\|^2 \geq 0$, $= 0$ iff $x = 0$

*Setup.* Let $f(x)$ be a function that maps from a Hilbert space to $\mathbb{R}$. We know that the derivative is the linear operator ("linear form") that takes a $dx$ (a infinitesimal change in the input) to $df$ (a infinitesimal change in the output):

$$df = f'(x)[dx].$$

*Riesz representation theorem* tells us that if we have a linear function that's "vector in number out", then it can be represented as a dot product with its input. So $f'(x)[dx]$ can be represented as the dot product between some vector and $dx$, and we call this vector the *gradient*:

$$df = f'(x)[dx] = (\nabla f) \cdot (dx).$$

An observation is that the gradient would always have the same "shape" as $x$.

*General strategy from deriving the gradient.* Start with $df$. Gradually massage it into a dot product between something and $dx$. That "something" would be the gradient.

Below we show some examples of this procedure from matrix-scalar functions. Note that the stuff on the LHS is equivalent to the stuff on the RHS, which is written in 18.06 style:

$$
\begin{aligned}
df &= (\nabla f) \cdot \text{dA} \\
&= \text{tr}((\nabla f)^T dA)
\end{aligned}
\quad \underset{\text{equivalent}}{\Leftrightarrow} \quad
\nabla f =
\begin{pmatrix}
\frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} & \cdots \\
\frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} & \cdots \\
\vdots & \vdots & \ddots
\end{pmatrix}
$$

# Example: gradient of $\|A\|_F$ with respect to $A$

Function:

$$f(A_{m \times n}) = \|A\|_F = \sqrt{\operatorname{tr}(A^T A)}$$

Deriving the gradient:

$$
\begin{aligned}
df &= \frac{1}{2\sqrt{\operatorname{tr}(A^T A)}} d(\operatorname{tr}(A^T A)) \quad \text{(scalar chain rule)} \\
&= \frac{1}{2\sqrt{\operatorname{tr}(A^T A)}} \operatorname{tr}[d(A^T A)] \quad \text{(trace is linear)} \\
&= \frac{1}{2\sqrt{\operatorname{tr}(A^T A)}} \operatorname{tr}[(dA)^T A + A^T dA] \quad \text{(matrix chain rule)} \\
&= \frac{1}{2\sqrt{\operatorname{tr}(A^T A)}} \Big( \operatorname{tr}[(dA)^T A] + \operatorname{tr}(A^T dA) \Big) \\
&= \frac{1}{2\sqrt{\operatorname{tr}(A^T A)}} \Big( \operatorname{tr}[A^T(dA)] + \operatorname{tr}(A^T dA) \Big) \quad (\operatorname{tr}(A) = \operatorname{tr}(A^T)) \\
&= \frac{1}{\sqrt{\operatorname{tr}(A^T A)}} \operatorname{tr}[A^T(dA)] \\
&= \frac{1}{\sqrt{\operatorname{tr}(A^T A)}} A \cdot dA \quad \text{(definition of the matrix dot product)} \\
&= \underbrace{\frac{A}{\sqrt{\operatorname{tr}(A^T A)}}}_{\text{gradient}} \cdot dA
\end{aligned}
$$

Note that the gradient is simply $A$ divided by $\|A\|_F$.

# Example: gradient of $x^T A y$ with respect to $A$

Function:

$$f(A_{m \times n}) = x^T A y$$

for some constant $x \in \mathbb{R}^m$ and $y \in \mathbb{R}^n$.

Deriving the gradient:

$$
\begin{aligned}
df &= x^T dA\, y \quad \text{(matrix product rule)} \\
&= \operatorname{tr}(x^T dA\, y) \\
&= \operatorname{tr}(y x^T dA) \\
&= \underbrace{(x y^T)}_{\text{gradient}} \cdot dA
\end{aligned}
$$

# Example: gradient of $\operatorname{sum}(A)$ with respect to $A$

Function:

$$f(A_{m \times n}) = \operatorname{sum}(A) = \mathbf{1}^T A \mathbf{1} = \operatorname{sum}(\operatorname{matrix}(1) \odot A)$$

Deriving the gradient (two ways are pretty much equivalent):

First way (using $\mathbf{1}^T A \mathbf{1}$):

$$\begin{aligned}
df &= d(\mathbf{1}^T A \mathbf{1}) \\
&= \mathbf{1}^T d A \mathbf{1} \\
&= \text{tr}(\mathbf{1}^T d A \mathbf{1}) \\
&= \text{tr}(\mathbf{1}\mathbf{1}^T d A) \\
&= \underbrace{(\mathbf{1}\mathbf{1}^T)}_{\text{gradient}} \cdot d A
\end{aligned}$$

Second way (using $\text{sum}(\text{matrix}(1) \odot A)$):

$$\begin{aligned}
df &= \text{sum}(\text{matrix}(1) \odot d A) \quad \text{(both sum and hadamard product are linear)} \\
&= \underbrace{\text{matrix}(1)}_{\text{gradient}} \cdot d A
\end{aligned}$$

# Lingering questions

- What would be the interpretation of the gradient if I define a weird but valid inner product?