

Matrix calculus

Ka Yu Wong

Two natural challenges naturally arise in calculating the derivative of a quantity with respect to a matrix. First, if the quantity is not a scalar, the result cannot be represented by a matrix. Second, the variable matrix has shape information (i.e. number of rows and columns) that we want to retain in the result. The use of differentials solves both problems elegantly by packing the derivative of a scalar with respect to a matrix (which is a matrix) into a differential form that captures all the derivative information.

Documents I found online are often incomplete, which motivates me to write a rigorous summary of the technique.

1 Differential form

Definition 1.1. \mathbb{R}^n is the n -dimensional Euclidean space. $e_i \in \mathbb{R}^n$ is the i -th standard basis vector of \mathbb{R}^n .

Definition 1.2. $(\mathbb{R}^n)^*$ is the dual space of \mathbb{R}^n , that is, the set of linear functions from \mathbb{R}^n to \mathbb{R} .

Definition 1.3. $x_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is the i -th coordinate function, that is, $x_i(e_j) = \delta_{ij}$.

Proposition 1.4. $(\mathbb{R}^n)^*$ is an n -dimensional vector space. $(x_i)_{i=1}^n$ is a basis of $(\mathbb{R}^n)^*$.

Definition 1.5. A 1-form on \mathbb{R}^n is a function $\omega : \mathbb{R}^n \rightarrow (\mathbb{R}^n)^*$.

Remark 1.6. A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and a 1-form ω on \mathbb{R}^n has a natural product as a

1-form: $f\omega : \mathbb{R}^n \rightarrow (\mathbb{R}^n)^*$, $f\omega(p) = f(p)\omega(p)$.

Definition 1.7. $dx_i : \mathbb{R}^n \rightarrow (\mathbb{R}^n)^*$ is the constant 1-form that sends all $p \in \mathbb{R}^n$ to x_i . That is, $dx_i(p) = x_i$ for all $p \in \mathbb{R}^n$.

Definition 1.8. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be smooth. Define $df : \mathbb{R}^n \rightarrow (\mathbb{R}^n)^*$,

$$df = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i.$$

Proposition 1.9. For every 1-form ω on \mathbb{R}^n , there exist unique $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$\omega = \sum_{i=1}^n f_i dx_i.$$

Proof. This follows directly from that for each $p \in \mathbb{R}^n$, $\omega(p)$ has a unique representation by $(x_i)_{i=1}^n$. □

This uniqueness is the key to why we can use differentials to fully capture the derivative.

2 Matrix calculus

Now we consider functions whose variables and values are matrices. We preserve those information in our definitions of differentials. Let $f, g : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times q}$, $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{q \times r}$.

Definition 2.1. Define $X : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ as the identity map:

$$X = \begin{bmatrix} x_{11} & \dots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \dots & x_{mn} \end{bmatrix}.$$

Denote f also as $[f]$:

$$f = [f] = \begin{bmatrix} f_{11} & \dots & f_{1q} \\ \vdots & \ddots & \vdots \\ f_{p1} & \dots & f_{pq} \end{bmatrix}.$$

Define $[f]^T : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{q \times p}$, $[f]^T(X) = f(X)^T$. Define $[f][h] : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times r}$, $[f][h](X) = f(X)h(X)$, the matrix product of images. If $p = q$ and $f(X)$ is invertible for all $X \in \mathbb{R}^{m \times n}$,

then define $[f]^{-1} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times p}$, $[f]^{-1}(X) = f(X)^{-1}$. If $p = q$, define $|f| : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, $|f|(X) = |f(X)|$, the determinant of its value.

Definition 2.2. Define df as a function from $\mathbb{R}^{m \times n}$ to a $p \times q$ matrix of 1-forms:

$$df = \begin{bmatrix} df_{11} & \dots & df_{1q} \\ \vdots & \ddots & \vdots \\ df_{p1} & \dots & df_{pq} \end{bmatrix}.$$

Specifically, dX is a function from $\mathbb{R}^{m \times n}$ to a $m \times n$ matrix of 1-forms:

$$dX = \begin{bmatrix} dx_{11} & \dots & dx_{1n} \\ \vdots & \ddots & \vdots \\ dx_{m1} & \dots & dx_{mn} \end{bmatrix}.$$

Definition 2.3. When $p = q = 1$, that is, $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, define the gradient $\frac{\partial f}{\partial X} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$:

$$\frac{\partial f}{\partial X} = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \dots & \frac{\partial f}{\partial x_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \dots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix}.$$

Proposition 2.4. Some basic properties. Let A, B be constant matrices.

- (a) $dA = 0$
- (b) $d(A[f]B) = A(df)B$
- (c) $d(f + g) = df + dg$
- (d) $d(f^T) = (df)^T$
- (e) $d(\text{tr}(f)) = \text{tr}(df)$
- (f) $d([f][h]) = (df)h + f(dh)$
- (g) $d[f]^{-1} = -[f]^{-1}(df)[f]^{-1}$

Proof. We only prove (f) and (g).

Proof for (f). If f, h are both scalars, then $d(fh) = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial fg}{\partial x_{ij}} dx_{ij} = \sum_{i=1}^m \sum_{j=1}^n (\frac{\partial f}{\partial x_{ij}} h + f \frac{\partial h}{\partial x_{ij}}) dx_{ij} = (df)h + f(dh)$. In general, let $1 \leq k \leq p, 1 \leq l \leq r$. Then we have $d([f][h])_{kl} = d(\sum_{u=1}^q f_{ku} h_{ul}) = \sum_{u=1}^q d(f_{ku} h_{ul}) = \sum_{u=1}^q (df_{ku}) h_{ul} + f_{ku} (dh_{ul}) = ((df)h)_{kl} + (f(dh))_{kl}$

Proof for (g). $0 = d(I_p) = d([f][f]^{-1}) = (df)[f]^{-1} + [f] d[f]^{-1}$. \square

We are ready to state the main result.

Theorem 2.5. Let $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, $g : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$. Then $g = \frac{\partial f}{\partial X}$ if and only if $df = \text{tr}(g^T dX)$.

Proof. If $A, B \in \mathbb{R}^{m \times n}$, then $\text{tr}(A^T B) = \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij}$, the sum of elementwise product. This is called the Hilbert–Schmidt inner product. Then $\text{tr}(g^T dX) = \sum_{i=1}^m \sum_{j=1}^n g_{ij} dx_{ij}$. By definition, $df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial x_{ij}} dx_{ij}$. The theorem then follows from Proposition 1.9. \square

This theorem is useful when we try to find the derivative of a scalar with respect to a matrix. We can write out the df and infer the form of the derivative.

Example 2.6. Let $a, b \in \mathbb{R}^m$ be constant. Compute the derivative of $a^T \Sigma^{-1} b$ with respect to Σ . We write out the differential form:

$$\begin{aligned} d(a^T \Sigma^{-1} b) &= d(\text{tr}(a^T \Sigma^{-1} b)) \\ &= \text{tr}(d(a^T \Sigma^{-1} b)) \\ &= \text{tr}(a^T (d\Sigma^{-1}) b) \\ &= \text{tr}(a^T (-\Sigma^{-1} (d\Sigma) \Sigma^{-1}) b) \\ &= \text{tr}(-\Sigma^{-1} b a^T \Sigma^{-1} d\Sigma) \end{aligned}$$

Hence $\frac{\partial a^T \Sigma^{-1} b}{\partial \Sigma} = (-\Sigma^{-1} b a^T \Sigma^{-1})^T = -\Sigma^{-T} a b^T \Sigma^{-T}$.

Proposition 2.7. Recall that X is the identity map, and $|X| : \mathbb{R}^{m \times m} \rightarrow \mathbb{R}$ sends a matrix to its determinant. We have $d|X| = |X| \text{tr}([X]^{-1} dX)$.

Proof. Let $1 \leq i \leq m, 1 \leq j \leq n$. The Laplace expansion along the i -th row gives $\frac{\partial |X|}{\partial x_{ij}} = (-1)^{i+j} |X_{-i,j}|$, where $X_{-i,j}$ is the matrix with i -th row and j -th column removed. This is precisely the i, j -th entry of the classical adjoint of X , so $\frac{\partial |X|}{\partial X} = |X|([X]^{-1})^T$. Then by Theorem 2.5, $d|X| = \text{tr}((|X|([X]^{-1})^T)^T dX) = |X| \text{tr}([X]^{-1} dX)$. \square

Now we present the chain rule.

Theorem 2.8. Let $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times q}, g : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}$. Then

$$d(g \circ f) = \text{tr}\left(\left(\frac{\partial g}{\partial f}\right)^T df\right).$$

Proof.

$$\begin{aligned} \left(\frac{\partial g}{\partial f}\right)^T df &= \sum_{k=1}^p \sum_{l=1}^q \frac{\partial g}{\partial f_{kl}} df_{kl} \\ &= \sum_{k=1}^p \sum_{l=1}^q \sum_{i=1}^m \sum_{j=1}^n \frac{\partial g}{\partial f_{kl}} \frac{\partial f_{kl}}{\partial x_{ij}} dx_{ij} \\ &= \sum_{i=1}^m \sum_{j=1}^n \left(\sum_{k=1}^p \sum_{l=1}^q \frac{\partial g}{\partial f_{kl}} \frac{\partial f_{kl}}{\partial x_{ij}} \right) dx_{ij} \\ &= \sum_{i=1}^m \sum_{j=1}^n \frac{\partial g \circ f}{\partial x_{ij}} dx_{ij} \\ &= \left(\frac{\partial g \circ f}{\partial x}\right)^T dX \end{aligned}$$

\square

Example 2.9. $d \ln |X| = \text{tr}\left(\frac{1}{|X|} d|X|\right) = \text{tr}\left(\frac{1}{|X|} |X| X^{-1} dX\right) = \text{tr}(X^{-1} dX)$. Hence $\frac{\partial \ln |X|}{\partial X} = X^{-T}$.

Example 2.10. Compute $\frac{\partial |X^T X|}{\partial X}$.

$$\begin{aligned}
d|X^T X| &= \text{tr}(|X^T X|(X^T X)^{-1} d(X^T X)) \\
&= \text{tr}(|X^T X|(X^T X)^{-1} (dX^T)X) + \text{tr}(|X^T X|(X^T X)^{-1} X^T dX) \\
&= \text{tr}((|X^T X|(X^T X)^{-1} (dX^T)X)^T) + \text{tr}(|X^T X|(X^T X)^{-1} X^T dX) \\
&= \text{tr}(X^T dX(|X^T X|(X^T X)^{-1}) + \text{tr}(|X^T X|(X^T X)^{-1} X^T dX) \\
&= 2|X^T X|\text{tr}((X(X^T X)^{-1})^T dX)
\end{aligned}$$

Hence $\frac{\partial |X^T X|}{\partial X} = 2|X^T X|X(X^T X)^{-1}$.

References

- [1] Kaare Brandt Petersen, Michael Syskind Pedersen, *The Matrix Cookbook*. <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>
- [2] Pili HU, *Matrix Calculus: Derivation and Simple Application*. <https://project.hupili.net/tutorial/hu2012-matrix-calculus/hu2012matrix-calculus.pdf>