

Matrix Calculus: Derivation and Simple Application

HU, Pili*

March 30, 2012[†]

Abstract

Matrix Calculus[3] is a very useful tool in many engineering problems. Basic rules of matrix calculus are nothing more than ordinary calculus rules covered in undergraduate courses. However, using matrix calculus, the derivation process is more compact. This document is adapted from the notes of a course the author recently attends. It builds matrix calculus from scratch. Only prerequisites are basic calculus notions and linear algebra operation. To get a quick executive guide, please refer to the cheat sheet in the end.

*hupili [at] ie [dot] cuhk [dot] edu [dot] hk

[†]Last compile: April 1, 2012

Contents

1	Introductory Example	3
2	Derivation	4
2.1	Organization of Elements	4
2.2	Deal with Inner Product	4
2.3	Properties of Trace	5
2.4	Deal with Generalized Inner Product	5
2.5	Define Matrix Differential	6
2.6	Matrix Differential Properties	8
2.7	Schema of Handling Scalar Function	9
2.8	Determinant	10
3	Application	11
3.1	The 2nd Induced Norm of Matrix	11
4	Cheat Sheet	12
	Acknowledgements	12
	References	12
	Appendix	13

1 Introductory Example

We start with an one variable linear function:

$$f(x) = ax \tag{1}$$

To be coherent, we abuse the partial derivative notation:

$$\frac{\partial f}{\partial x} = a \tag{2}$$

Extending this function to be multivariate, we have:

$$f(x) = \sum_i a_i x_i = a^T x \tag{3}$$

Where $a = [a_1, a_2, \dots, a_n]^T$ and $x = [x_1, x_2, \dots, x_n]^T$. We first compute partial derivatives directly:

$$\frac{\partial f}{\partial x_k} = \frac{\partial(\sum_i a_i x_i)}{\partial x_k} = a_k \tag{4}$$

for all $k = 1, 2, \dots, n$. Then we organize n partial derivatives in the following way:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = a \tag{5}$$

The first equality is by proper definition and the rest roots from ordinary calculus rules.

Eqn(5) is analogous to eqn(2), except the variable changes from a scalar to a vector. Thus we want to directly claim the result of eqn(5) without those intermediate steps solving for partial derivatives separately. Actually, we'll see soon that eqn(5) plays a core role in matrix calculus.

Following sections are organized as follows:

- Section(2) builds commonly used matrix calculus rules from ordinary calculus and linear algebra. Necessary and important properties of linear algebra is also proved along the way. This section is not organized afterward. All results are proved when we need them.
- Section(3) shows some applications using matrix calculus.
- Section(4) concludes a cheat sheet of matrix calculus. Note that this cheat sheet may be different from others. Users need to figure out some basic definitions before applying the rules.

2 Derivation

2.1 Organization of Elements

From the introductory example, we already see that matrix calculus does not distinguish from ordinary calculus by fundamental rules. However, with better organization of elements and proving useful properties, we can simplify the derivation process in real problems.

The author would like to adopt the following definition:

Definition 1. For a scalar valued function $f(x)$, the result $\frac{\partial f}{\partial x}$ has the same size with x . That is

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{12}} & \cdots & \frac{\partial f}{\partial x_{1n}} \\ \frac{\partial f}{\partial x_{21}} & \frac{\partial f}{\partial x_{22}} & \cdots & \frac{\partial f}{\partial x_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \frac{\partial f}{\partial x_{m2}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix} \quad (6)$$

In eqn(2), x is a 1-by-1 matrix and the result $\frac{\partial f}{\partial x} = a$ is also a 1-by-1 matrix. In eqn(5), x is a column vector (known as n-by-1 matrix) and the result $\frac{\partial f}{\partial x} = a$ has the same size.

Example 1. By this definition, we have:

$$\frac{\partial f}{\partial x^T} = \left(\frac{\partial f}{\partial x}\right)^T = a^T \quad (7)$$

Note that we only use the organization definition in this example. Later we'll show that with some matrix properties, this formula can be derived without using $\frac{\partial f}{\partial x}$ as a bridge.

2.2 Deal with Inner Product

Theorem 1. If there's a multivariate scalar function $f(x) = a^T x$, we have $\frac{\partial f}{\partial x} = a$.

Proof. See introductory example. \square

Since $a^T x$ is scalar, we can write it equivalently as the trace of its own. Thus,

Proposition 2. If there's a multivariate scalar function $f(x) = \text{Tr}[a^T x]$, we have $\frac{\partial f}{\partial x} = a$.

$\text{Tr}[\bullet]$ is the operator to sum up diagonal elements of a matrix. In the next section, we'll explore more properties of trace. As long as we can transform our target function into the form of theorem(1) or proposition(2), the result can be written out directly. Notice in proposition(2), a and x are both vectors. We'll show later as long as their sizes agree, it holds for matrix a and x .

2.3 Properties of Trace

Definition 2. Trace of square matrix is defined as: $\text{Tr}[A] = \sum_i A_{ii}$

Theorem 3. Matrix trace has the following properties:

- (1) $\text{Tr}[A + B] = \text{Tr}[A] + \text{Tr}[B]$
- (2) $\text{Tr}[cA] = c\text{Tr}[A]$
- (3) $\text{Tr}[AB] = \text{Tr}[BA]$
- (4) $\text{Tr}[A_1 A_2 \dots A_n] = \text{Tr}[A_n A_1 \dots A_{n-1}]$
- (5) $\text{Tr}[A^T B] = \sum_i \sum_j A_{ij} B_{ij}$
- (6) $\text{Tr}[A] = \text{Tr}[A^T]$

where A, B are matrices with proper sizes, and c is a scalar value.

Proof. See wikipedia [5] for the proof. □

Here we explain the intuitions behind each property to make it easier to remember. Property(1) and property(2) shows the linearity of trace. Property(3) means two matrices' multiplication inside a the trace operator is commutative. Note that the matrix multiplication without trace is not commutative and the commutative property inside the trace does not hold for more than 2 matrices. Property (4) is the proposition of property (3) by considering $A_1 A_2 \dots A_{n-1}$ as a whole. It is known as cyclic property, so that you can rotate the matrices inside a trace operator. Property (5) shows a way to express the sum of element by element product using matrix product and trace. Note that inner product of two vectors is also the sum of element by element product. Property (5) resembles the vector inner product by $\text{form}(A^T B)$. The author regards property (5) as the extension of inner product to matrices(Generalized Inner Product).

2.4 Deal with Generalized Inner Product

Theorem 4. If there's a multivariate scalar function $f(x) = \text{Tr}[A^T x]$, we have $\frac{\partial f}{\partial x} = A$. (A, x can be matrices).

Proof. Using property (5) of trace, we can write f as:

$$f(x) = \text{Tr}[A^T x] = \sum_{ij} A_{ij} x_{ij} \quad (8)$$

It's easy to show:

$$\frac{\partial f}{\partial x_{ij}} = \frac{\partial(\sum_{ij} A_{ij} x_{ij})}{\partial x_{ij}} = A_{ij} \quad (9)$$

Organize elements using definition(1), it is proved. □

With this theorem and properties of trace we revisit example(1).

Example 2. For vector a, x and function $f(x) = a^T x$

$$\frac{\partial f}{\partial x^T} \quad (10)$$

$$= \frac{\partial(a^T x)}{\partial x^T} \quad (11)$$

$$(f \text{ is scalar}) = \frac{\partial(\text{Tr}[a^T x])}{\partial x^T} \quad (12)$$

$$(property(3)) = \frac{\partial(\text{Tr}[x a^T])}{\partial x^T} \quad (13)$$

$$(property(6)) = \frac{\partial(\text{Tr}[a x^T])}{\partial x^T} \quad (14)$$

$$(property \text{ of transpose}) = \frac{\partial(\text{Tr}[(a^T)^T x^T])}{\partial x^T} \quad (15)$$

$$(theorem(4)) = a^T \quad (16)$$

The result is the same with example(1), where we used the basic definition.

The above example actually demonstrates the usual way of handling a matrix derivative problem.

2.5 Define Matrix Differential

Although we want matrix derivative at most time, it turns out matrix differential is easier to operate due to the form invariance property of differential. Matrix differential inherit this property as a natural consequence of the following definition.

Definition 3. Define matrix differential:

$$dA = \begin{bmatrix} dA_{11} & dA_{12} & \dots & dA_{1n} \\ dA_{21} & dA_{22} & \dots & dA_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ dA_{m1} & dA_{m2} & \dots & dA_{mn} \end{bmatrix} \quad (17)$$

Theorem 5. Differential operator is distributive through trace operator:
 $d\text{Tr}[A] = \text{Tr}[dA]$

Proof.

$$\text{LHS} = d\left(\sum_i A_{ii}\right) = \sum_i dA_{ii} \quad (18)$$

$$\text{RHS} = \text{Tr} \begin{bmatrix} dA_{11} & dA_{12} & \dots & dA_{1n} \\ dA_{21} & dA_{22} & \dots & dA_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ dA_{m1} & dA_{m2} & \dots & dA_{mn} \end{bmatrix} \quad (19)$$

$$= \sum_i dA_{ii} = \text{LHS} \quad (20)$$

□

Now that matrix differential is well defined, we want to relate it back to matrix derivative. The scalar version differential and derivative can be related as follows:

$$df = \frac{\partial f}{\partial x} dx \quad (21)$$

So far, we're dealing with scalar function f and matrix variable x . $\frac{\partial f}{\partial x}$ and dx are both matrix according to definition. In order to make the quantities in eqn(21) equal, we must figure out a way to make the RHS a scalar. It's not surprising that trace is what we want.

Theorem 6.

$$df = \text{Tr} \left[\left(\frac{\partial f}{\partial x} \right)^T dx \right] \quad (22)$$

for scalar function f and arbitrarily sized x .

Proof.

$$\text{LHS} = df \quad (23)$$

$$\text{(definition of scalar differential)} = \sum_{ij} \frac{\partial f}{\partial x_{ij}} dx_{ij} \quad (24)$$

$$\text{RHS} = \text{Tr} \left[\left(\frac{\partial f}{\partial x} \right)^T dx \right] \quad (25)$$

$$\text{(trace property (5))} = \sum_{ij} \left(\frac{\partial f}{\partial x} \right)_{ij} (dx)_{ij} \quad (26)$$

$$\text{(definition(3))} = \sum_{ij} \left(\frac{\partial f}{\partial x} \right)_{ij} dx_{ij} \quad (27)$$

$$\text{(definition(1))} = \sum_{ij} \frac{\partial f}{\partial x_{ij}} dx_{ij} \quad (28)$$

$$= \text{LHS} \quad (29)$$

□

Theorem(6) is the bridge between matrix derivative and matrix differential. We'll see in later applications that matrix differential is more convenient to manipulate. After certain manipulation we can get the form of theorem(6). Then we can directly write out matrix derivative using this theorem.

2.6 Matrix Differential Properties

Theorem 7. *We claim the following properties of matrix differential:*

- $d(cA) = cdA$
- $d(A + B) = dA + dB$
- $d(AB) = dAB + AdB$

Proof. They're all natural consequences given the definition(3). We only show the 3rd one in this document. Note that the equivalence holds if LHS and RHS are equivalent element by element. We consider the (ij)-th element.

$$\text{LHS}_{ij} = d\left(\sum_k A_{ik}B_{kj}\right) \quad (30)$$

$$= \sum_k (dA_{ik}B_{kj} + A_{ik}dB_{kj}) \quad (31)$$

$$\text{RHS}_{ij} = (dAB)_{ij} + (AdB)_{ij} \quad (32)$$

$$= \sum_k dA_{ik}B_{kj} + \sum_k A_{ik}dB_{kj} \quad (33)$$

$$= \text{LHS}_{ij} \quad (34)$$

□

Example 3. *Given the function $f(x) = x^T Ax$, where A is square and x is a column vector, we can compute:*

$$df = d\text{Tr}[x^T Ax] \quad (35)$$

$$= \text{Tr}[d(x^T Ax)] \quad (36)$$

$$= \text{Tr}[d(x^T)Ax + x^T d(Ax)] \quad (37)$$

$$= \text{Tr}[d(x^T)Ax + x^T dAx + x^T Adx] \quad (38)$$

$$(A \text{ is constant}) = \text{Tr}[dx^T Ax + x^T Adx] \quad (39)$$

$$= \text{Tr}[dx^T Ax] + \text{Tr}[x^T Adx] \quad (40)$$

$$= \text{Tr}[x^T A^T dx] + \text{Tr}[x^T Adx] \quad (41)$$

$$= \text{Tr}[x^T A^T dx + x^T Adx] \quad (42)$$

$$= \text{Tr}[(x^T A^T + x^T A)dx] \quad (43)$$

Using theorem(6), we obtain the derivative:

$$\frac{\partial f}{\partial x} = (x^T A^T + x^T A)^T = Ax + A^T x \quad (44)$$

When A is symmetric, it simplifies to:

$$\frac{\partial f}{\partial x} = 2Ax \quad (45)$$

Let $A = I$, we have:

$$\frac{\partial(x^T x)}{\partial x} = 2x \quad (46)$$

Example 4. For a non-singular square matrix X , we have $XX^{-1} = I$. Take matrix differentials at both sides:

$$0 = dI = d(XX^{-1}) = dXX^{-1} + Xd(X^{-1}) \quad (47)$$

Rearrange terms:

$$d(X^{-1}) = -X^{-1}dXX^{-1} \quad (48)$$

2.7 Schema of Handling Scalar Function

The above example already demonstrates the general schema. Here we conclude the process:

1. $df = d\text{Tr}[f] = \text{Tr}[df]$
2. Apply trace properties(see theorem(3)) and matrix differential properties(see theorem(7)) to get the following form:

$$df = \text{Tr}[A^T x] \quad (49)$$

3. Apply theorem(6) to get:

$$\frac{\partial f}{\partial x} = A \quad (50)$$

To this point, you can handle many problems. In this schema, matrix differential and trace play crucial roles. Later we'll deduce some widely used formula to facilitate potential applications. As you will see, although we rely on matrix differential in the schema, the deduction of certain formula may be more easily done using matrix derivatives.

2.8 Determinant

For a background of determinant, please refer to [6]. We first quote some definitions and properties without proof:

Theorem 8. *Let A be a square matrix:*

- *The minor M_{ij} is obtained by remove i -th row and j -th column of A and then take determinant of the resulting $(n-1)$ by $(n-1)$ matrix.*
- *The ij -th cofactor is defined as $C_{ij} = (-1)^{i+j} M_{ij}$.*
- *If we expand determinant with respect to the i -th row, $\det(A) = \sum_j A_{ij} C_{ij}$.*
- *The adjugate of A is defined as $\text{adj}(A)_{ij} = (-1)^{i+j} M_{ji} = C_{ji}$. So $\text{adj}(A) = C^T$*
- *For non-singular matrix A , we have: $A^{-1} = \frac{\text{adj}(A)}{\det(A)} = \frac{C^T}{\det(A)}$*

Now we're ready to show the derivative of determinant. Note that determinant is just a scalar function, so all techniques discussed above is applicable. We first write the derivative element by element. Expanding determinant on the i -th row, we have:

$$\frac{\partial \det(A)}{\partial A_{ij}} = \frac{\partial (\sum_j A_{ij} C_{ij})}{\partial A_{ij}} = C_{ij} \quad (51)$$

First equality is from determinant definition and second equality is by the observation that only coefficient of A_{ij} is left. Grouping all elements using definition(1), we have:

$$\frac{\partial \det(A)}{\partial A} = C = \text{adj}(A)^T \quad (52)$$

If A is non-singular, we have:

$$\frac{\partial \det(A)}{\partial A} = (\det(A) A^{-1})^T = \det(A) (A^{-1})^T \quad (53)$$

Next, we use theorem(6) to give the differential relationship:

$$d \det(A) = \text{Tr} \left[\left(\frac{\partial \det(A)}{\partial A} \right)^T dA \right] \quad (54)$$

$$= \text{Tr} [(\det(A) (A^{-1})^T)^T dA] \quad (55)$$

$$= \text{Tr} [\det(A) A^{-1} dA] \quad (56)$$

In many practical problem, the log determinant is more widely used:

$$\frac{\partial \ln \det(A)}{\partial A} = \frac{1}{\det(A)} \frac{\partial \det(A)}{\partial A} = (A^{-1})^T \quad (57)$$

The first equality comes from chain rule of ordinary calculus($\ln \det(A)$ and $\det(A)$ are both scalars). Similarly, we derive for differential:

$$d \ln \det(A) = \text{Tr} [A^{-1} dA] \quad (58)$$

3 Application

3.1 The 2nd Induced Norm of Matrix

The induced norm of matrix is defined as [4]:

$$\|A\|_p = \max_x \frac{\|Ax\|_p}{\|x\|_p} \quad (59)$$

where $\|\bullet\|_p$ denotes the p-norm of vectors. Now we solve for $p = 2$. (By default, $\|\bullet\|$ means $\|\bullet\|_2$)

The problem can be restated as:

$$\|A\|^2 = \max_x \frac{\|Ax\|^2}{\|x\|^2} \quad (60)$$

since all quantities involved are non-negative. Then we consider a scaling of vector $x' = tx$, thus:

$$\|A\|^2 = \max_{x'} \frac{\|Ax'\|^2}{\|x'\|^2} = \max_x \frac{\|tAx\|^2}{\|tx\|^2} = \max_x \frac{t^2\|Ax\|^2}{t^2\|x\|^2} = \max_x \frac{\|Ax\|^2}{\|x\|^2} \quad (61)$$

This shows the invariance under scaling. Now we can restrict our attention to those x with $\|x\| = 1$, and reach the following formulation:

$$\text{Maximize} \quad f(x) = \|Ax\|^2 \quad (62)$$

$$s.t. \quad \|x\|^2 = 1 \quad (63)$$

The standard way to handle this constrained optimization is using Lagrange relaxation:

$$L(x) = f(x) - \lambda(\|x\|^2 - 1) \quad (64)$$

Then we apply the general schema of handling scalar function on $L(x)$. First take differential:

$$dL(x) = d\text{Tr}[L(x)] \quad (65)$$

$$= \text{Tr}[d(L(x))] \quad (66)$$

$$= \text{Tr}[d(x^T A^T Ax - \lambda(x^T x - 1))] \quad (67)$$

$$= \text{Tr}[2x^T A^T A dx - \lambda(2x^T dx)] \quad (68)$$

$$= \text{Tr}[(2x^T A^T A - 2\lambda x^T)dx] \quad (69)$$

Next write out derivative:

$$\frac{\partial L}{\partial x} = 2A^T Ax - 2\lambda x \quad (70)$$

Let $\frac{\partial L}{\partial x} = 0$, we have:

$$(A^T A)x = \lambda x \quad (71)$$

That means x is the eigen vector of $(A^T A)$ (normalized to $\|x\| = 1$), and λ is corresponding eigen value. We plug this result back to objective function:

$$f(x) = x^T(A^T A x) = x^T(\lambda x) = \lambda \quad (72)$$

which means, to maximize $f(x)$, we should pick the maximum eigen value:

$$\|A\|^2 = \max_x f(x) = \lambda_{\max}(A^T A) \quad (73)$$

That is:

$$\|A\| = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max}(A) \quad (74)$$

where σ_{\max} denotes the maximum singular value. If A is real symmetric, $\sigma_{\max}(A) = \lambda_{\max}(A)$.

Now we consider a real symmetric A and check whether:

$$\lambda_{\max}^2(A) = \max_x \frac{\|Ax\|^2}{\|x\|^2} = \max_x \frac{x^T A^T A x}{x^T x} \quad (75)$$

Proof. Since $A^T A$ is real symmetric, it has an orthonormal basis formed by n eigen vectors, v_1, v_2, \dots, v_n , with eigen values $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. We can write $x = \sum_i c_i v_i$, where $c_i = \langle x, v_i \rangle$. Then,

$$\frac{x^T A^T A x}{x^T x} \quad (76)$$

$$(v_k \text{ is an orthonormal set}) = \frac{\sum_i \lambda_i c_i^2}{\sum_i c_i^2} \quad (77)$$

$$\leq \frac{\sum_i \lambda_1 c_i^2}{\sum_i c_i^2} \quad (78)$$

$$= \lambda_1 \quad (79)$$

Now we have proved an upper bound for $\|A\|^2$. We show this bound is achievable by assigning $x = v_1$. \square

4 Cheat Sheet

Acknowledgements

Thanks prof. XU, Lei's tutorial on matrix calculus. Besides, the author also benefit a lot from other online materials.

References

- [1] Pili Hu. Matrix calculus. GitHub, <https://github.com/hupili/tutorial/tree/master/matrix-calculus>, 3 2012. HU, Pili's tutorial collection.

- [2] Pili Hu. Tutorial collection. GitHub, <https://github.com/hupili/tutorial>, 3 2012. HU, Pili's tutorial collection.
- [3] Matrix Calculus, Wikipedia, http://en.wikipedia.org/wiki/Matrix_calculus
- [4] Matrix Norm, Wikipedia, http://en.wikipedia.org/wiki/Matrix_norm
- [5] Matrix Trace, Wikipedia, http://en.wikipedia.org/wiki/Trace_%28linear_algebra%29
- [6] Matrix Determinant, Wikipedia, <http://en.wikipedia.org/wiki/Determinant>

Appendix