

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 3, 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	6
2.5	Define Matrix Differential	7
2.6	Matrix Differential Properties	8
2.7	Schema of Handling Scalar Function	9
2.8	Determinant	10
2.9	Vector Function and Vector Variable	11
2.10	Vector Function Differential	13
2.11	Chain Rule	14
3	Application	16
3.1	The 2nd Induced Norm of Matrix	16
3.2	General Multivariate Gaussian Distribution	17
3.3	Maximum Likelihood Estimation of Gaussian	19
4	Cheat Sheet	21
	Acknowledgements	21
	References	21
	Appendix	22

1 Introductory Example

We start with an one variable linear function:

$$f(x) = ax \quad (1)$$

To be coherent, we abuse the partial derivative notation:

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

Extending this function to be multivariate, we have:

$$f(x) = \sum_i a_i x_i = a^T x \quad (3)$$

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:

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

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

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

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 afterhand. All results are proved when we need them.
- Section(3) shows some applications using matrix calculus. Table(1) shows the relation between Section(2) and Section(3).
- 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.

Table 1: Derivation and Application Correspondance

Derivation	Application
2.1-2.7	3.1
2.9,2.10	3.2
2.8,2.11	3.3

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: $\text{Tr}[A] = \sum_i A_{ii}$*

Example 2. *Using definition(1,2), it is very easy to show:*

$$\frac{\partial \text{Tr}[A]}{\partial A} = I \quad (8)$$

since only diagonal elements are kept by the trace operator.

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. \square

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 form $\text{form}(A^T B)$. 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:

$$f(x) = \text{Tr}[A^T x] = \sum_{ij} A_{ij} x_{ij} \quad (9)$$

It's easy to show:

$$\frac{\partial f}{\partial x_{ij}} = \frac{\partial(\sum_{ij} A_{ij} x_{ij})}{\partial x_{ij}} = A_{ij} \quad (10)$$

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

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

Example 3. *For vector a, x and function $f(x) = a^T x$*

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

$$= \frac{\partial(a^T x)}{\partial x^T} \quad (12)$$

$$(f \text{ is scalar}) = \frac{\partial(\text{Tr}[a^T x])}{\partial x^T} \quad (13)$$

$$(property(3)) = \frac{\partial(\text{Tr}[x a^T])}{\partial x^T} \quad (14)$$

$$(property(6)) = \frac{\partial(\text{Tr}[a x^T])}{\partial x^T} \quad (15)$$

$$(property \text{ of transpose}) = \frac{\partial(\text{Tr}[(a^T)^T x^T])}{\partial x^T} \quad (16)$$

$$(theorem(4)) = a^T \quad (17)$$

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:*

$$dA = \begin{bmatrix} dA_{11} & dA_{12} & \dots & dA_{1n} \\ dA_{21} & dA_{22} & \dots & dA_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ dA_{m1} & dA_{m2} & \dots & dA_{mn} \end{bmatrix} \quad (18)$$

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

Proof.

$$\text{LHS} = d\left(\sum_i A_{ii}\right) = \sum_i dA_{ii} \quad (19)$$

$$\text{RHS} = \text{Tr} \begin{bmatrix} dA_{11} & dA_{12} & \dots & dA_{1n} \\ dA_{21} & dA_{22} & \dots & dA_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ dA_{m1} & dA_{m2} & \dots & dA_{mn} \end{bmatrix} \quad (20)$$

$$= \sum_i dA_{ii} = \text{LHS} \quad (21)$$

□

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:

$$df = \frac{\partial f}{\partial x} dx \quad (22)$$

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(22) 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.

$$df = \text{Tr} \left[\left(\frac{\partial f}{\partial x} \right)^T dx \right] \quad (23)$$

for scalar function f and arbitrarily sized x .

Proof.

$$\text{LHS} = df \quad (24)$$

$$\text{(definition of scalar differential)} = \sum_{ij} \frac{\partial f}{\partial x_{ij}} dx_{ij} \quad (25)$$

$$\text{RHS} = \text{Tr} \left[\left(\frac{\partial f}{\partial x} \right)^T dx \right] \quad (26)$$

$$\text{(trace property (5))} = \sum_{ij} \left(\frac{\partial f}{\partial x} \right)_{ij} (dx)_{ij} \quad (27)$$

$$\text{(definition(3))} = \sum_{ij} \left(\frac{\partial f}{\partial x} \right)_{ij} dx_{ij} \quad (28)$$

$$\text{(definition(1))} = \sum_{ij} \frac{\partial f}{\partial x_{ij}} dx_{ij} \quad (29)$$

$$= \text{LHS} \quad (30)$$

□

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.

$$\text{LHS}_{ij} = d\left(\sum_k A_{ik} B_{kj}\right) \quad (31)$$

$$= \sum_k (dA_{ik} B_{kj} + A_{ik} dB_{kj}) \quad (32)$$

$$\text{RHS}_{ij} = (dAB)_{ij} + (AdB)_{ij} \quad (33)$$

$$= \sum_k dA_{ik} B_{kj} + \sum_k A_{ik} dB_{kj} \quad (34)$$

$$= \text{LHS}_{ij} \quad (35)$$

□

Example 4. Given the function $f(x) = x^T Ax$, where A is square and x is a column vector, we can compute:

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

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

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

$$= \text{Tr}[d(x^T)Ax + x^T dAx + x^T Adx] \quad (39)$$

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

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

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

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

$$= \text{Tr}[(x^T A^T + x^T A)dx] \quad (44)$$

Using theorem(6), we obtain the derivative:

$$\frac{\partial f}{\partial x} = (x^T A^T + x^T A)^T = Ax + A^T x \quad (45)$$

When A is symmetric, it simplifies to:

$$\frac{\partial f}{\partial x} = 2Ax \quad (46)$$

Let $A = I$, we have:

$$\frac{\partial(x^T x)}{\partial x} = 2x \quad (47)$$

Example 5. For a non-singular square matrix X , we have $XX^{-1} = I$. Take matrix differentials at both sides:

$$0 = dI = d(XX^{-1}) = dXX^{-1} + Xd(X^{-1}) \quad (48)$$

Rearrange terms:

$$d(X^{-1}) = -X^{-1}dXX^{-1} \quad (49)$$

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:

$$df = \text{Tr}[A^T x] \quad (50)$$

3. Apply theorem(6) to get:

$$\frac{\partial f}{\partial x} = A \quad (51)$$

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:

$$\frac{\partial \det(A)}{\partial A_{ij}} = \frac{\partial (\sum_j A_{ij} C_{ij})}{\partial A_{ij}} = C_{ij} \quad (52)$$

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:

$$\frac{\partial \det(A)}{\partial A} = C = \text{adj}(A)^T \quad (53)$$

If A is non-singular, we have:

$$\frac{\partial \det(A)}{\partial A} = (\det(A) A^{-1})^T = \det(A) (A^{-1})^T \quad (54)$$

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

$$d \det(A) = \text{Tr} \left[\left(\frac{\partial \det(A)}{\partial A} \right)^T dA \right] \quad (55)$$

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

$$= \text{Tr} [\det(A)A^{-1}dA] \quad (57)$$

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

$$\frac{\partial \ln \det(A)}{\partial A} = \frac{1}{\det(A)} \frac{\partial \det(A)}{\partial A} = (A^{-1})^T \quad (58)$$

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

$$d \ln \det(A) = \text{Tr} [A^{-1}dA] \quad (59)$$

2.9 Vector Function and Vector Variable

The above sections show how to deal with scalar functions. In order to deal with vector function, we should restrict our attention to vector variables. It's no surprising that the tractable forms in matrix calculus is so scarce. If we allow matrix functions and matrix variables, given the fact that fully specification of all partial derivatives calls for a tensor, it will be difficult to visualize the result on a paper. An alternative is to stretch functions and variables such that they appear as vectors.

An annoying fact of matrix calculus is that, when you try to find reference materials, there are always two kinds of people. One group calculates as the transpose of another. Many online resources are not coherent, which mislead people.

We borrow the following definitions of Hessian matrix[7] and Jacobian matrix[8] from Wikipedia:

$$H(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad (60)$$

$$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \quad (61)$$

Note two things about second order derivative:

- By writting the abbreviation $\frac{\partial^2 f}{\partial x_2 \partial x_1}$, people mean $\frac{\partial}{\partial x_2} \left(\frac{\partial f}{\partial x_1} \right)$ by convention. That is, first take derivative with respect to x_1 and then take derivative with respect to x_2 .
- The Hessian matrix can be regarded as first compute the $\frac{\partial f}{\partial x^T}$ (using definition(1) to organize), and then compute a vector-function-to-vector-variable derivative treating $\frac{\partial f}{\partial x^T}$ as the function.

Bearing this in mind, we find Hessian matrix and Jacobian matrix actually have contradictory notion of organization. In order to be coherent in this document, we adopt the Hessian style. That is, each row corresponds to a variable, and each column corresponds to a function. To be concrete:

Definition 4. For a vector function $f = [f_1, f_2, \dots, f_n]^T$, and $f_i = f_i(x)$ where $x = [x_1, x_2, \dots, x_m]^T$, we have the following definition:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_2}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_1} \\ \frac{\partial f_1}{\partial x_2} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_m} & \frac{\partial f_2}{\partial x_m} & \cdots & \frac{\partial f_n}{\partial x_m} \end{bmatrix} \quad (62)$$

Example 6. According to definition(4), we revisit the definition of Hessian and Jacobian.

Given twice differentiable function $f(x)$, the Hessian is defined as:

$$H(f) = \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial x^T} \right) \quad (63)$$

Given two variables x and y , if we want to transform x into y , the Jacobian is defined as:

$$J = \det \left(\left(\frac{\partial x}{\partial y} \right)^T \right) = \det \left(\frac{\partial x}{\partial y} \right) \quad (64)$$

Note "Jacobian" is the shorthand name of "Jacobian determinant", which is the determinant of "Jacobian matrix". Due to the transpose invariance of determinant, the second equality shows that it does not matter which organization method we use if we only want to do compute the Jacobian, rather than Jacobin matrix. However, if we're to write out the Jacobin matrix, this may be a pitfall depending on what organization of vector-to-vector derivative we define.

2.10 Vector Function Differential

In the previous sections, we relate matrix differential with matrix derivative using theorem(6). This theorem bridges the two quantities using trace. Thus we can handle the problem using preferred form interchangeably. (As is seen: calculating derivative of determinant is more direct; calculating differential of inverse is tractable.)

In this section, we provide another theorem to relate vector function differential to vector-to-vector derivatives. Amazingly, it takes a cleaner form.

Theorem 9. *Consider the definition(4), we have: $df = (\frac{\partial f}{\partial x})^T dx$*

Proof. Apparently, df has n components, so we prove element by element. Consider the j -th component:

$$\text{LHS}_j = df_j \tag{65}$$

$$= \sum_{i=1}^m \frac{\partial f_j}{\partial x_i} dx_i \tag{66}$$

$$\text{RHS}_j = ((\frac{\partial f}{\partial x})^T x)_j \tag{67}$$

$$= \sum_{i=1}^m (\frac{\partial f}{\partial x})_{ji}^T x_i \tag{68}$$

$$= \sum_{i=1}^m (\frac{\partial f}{\partial x})_{ij} x_i \tag{69}$$

$$= \sum_{i=1}^m (\frac{\partial f}{\partial x})_{ij} x_i \tag{70}$$

$$= \sum_{i=1}^m (\frac{\partial f_j}{\partial x_i}) x_i = \text{LHS}_j \tag{71}$$

□

Note that the trace operator is gone compared with theorem(6) due to the nice way of defining matrix vector multiplication. We can have a similar schema of handling vector-to-vector derivatives using this scheme. We don't bother to list the schema again. Instead, we provide an example.

Example 7. *Consider the variable transformation: $x = \sigma \Lambda^{-0.5} W^T \xi$, where σ is a real value, Λ is full rank diagonal matrix, and W is orthonormal square matrix (namely $WW^T = W^T W = I$). Compute the absolute value of Jacobian.*

First, we want to find $\frac{\partial x}{\partial \xi}$. This can be easily done by computing the differential:

$$dx = d(\sigma \Lambda^{-0.5} W^T \xi) = \sigma \Lambda^{-0.5} W^T d\xi \quad (72)$$

Applying theorem(9), we have:

$$\frac{\partial x}{\partial \xi} = (\sigma \Lambda^{-0.5} W^T)^T \quad (73)$$

Thus,

$$J_m = \left(\frac{\partial x}{\partial \xi}\right)^T = ((\sigma \Lambda^{-0.5} W^T)^T)^T = \sigma \Lambda^{-0.5} W^T \quad (74)$$

where J_m is the Jacobian matrix and $J = \det(J_m)$. Then we use some property of determinant to calculate the absolute value of Jacobian:

$$|J| = |\det(J_m)| \quad (75)$$

$$= \sqrt{|\det(J_m)| |\det(J_m)|} \quad (76)$$

$$= \sqrt{|\det(J_m^T)| |\det(J_m)|} \quad (77)$$

$$= \sqrt{|\det(J_m^T J_m)|} \quad (78)$$

$$= \sqrt{|\det(J_m J_m^T)|} \quad (79)$$

$$= \sqrt{|\det(W \Lambda^{-0.5} \sigma \sigma \Lambda^{-0.5} W^T)|} \quad (80)$$

$$= \sqrt{|\det(\sigma^2 W \Lambda^{-1} W^T)|} \quad (81)$$

If the dimension of x and ξ is d and we define $\Sigma = W \Lambda W^T$. A nice result is calculated:

$$|J| = \sigma^d \det(\Sigma)^{-1/2} \quad (82)$$

which we'll see an application of generalizing the multivariate Gaussian distribution.

2.11 Chain Rule

In the schemas we conclude above, differential is convenient in many problems. For derivative, the nice aspect is that we have chain rule, which is an analogy to the chain rule in ordinary calculus. However, one should be very careful applying this chain rule, for the multiplication of matrix requires dimension agreement.

Theorem 10. Suppose we have n column vectors $x^{(1)}, x^{(2)}, \dots, x^{(n)}$, each is of length l_1, l_2, \dots, l_n . We know $x^{(i)}$ is a function of $x^{(i-1)}$, for all $i = 2, 3, \dots, n$. The following relationship holds:

$$\frac{\partial x^{(n)}}{\partial x^{(1)}} = \frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \cdots \frac{\partial x^{(n)}}{\partial x^{(n-1)}} \quad (83)$$

Proof. Under definition(4), theorem(9) holds. We apply this theorem to each consecutive vectors:

$$\mathrm{d}x^{(2)} = \left(\frac{\partial x^{(2)}}{\partial x^{(1)}}\right)^T \mathrm{d}x^{(1)} \quad (84)$$

$$\mathrm{d}x^{(3)} = \left(\frac{\partial x^{(3)}}{\partial x^{(2)}}\right)^T \mathrm{d}x^{(2)} \quad (85)$$

$$\dots \quad (86)$$

$$\mathrm{d}x^{(n)} = \left(\frac{\partial x^{(n)}}{\partial x^{(n-1)}}\right)^T \mathrm{d}x^{(n-1)} \quad (87)$$

Plug previous one into next one, we have:

$$\mathrm{d}x^{(n)} = \left(\frac{\partial x^{(n)}}{\partial x^{(n-1)}}\right)^T \dots \left(\frac{\partial x^{(3)}}{\partial x^{(2)}}\right)^T \left(\frac{\partial x^{(2)}}{\partial x^{(1)}}\right)^T \mathrm{d}x^{(1)} \quad (88)$$

$$= \left(\frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \dots \frac{\partial x^{(n)}}{\partial x^{(n-1)}}\right)^T \mathrm{d}x^{(1)} \quad (89)$$

Applying theorem(9) again in the reverse direction, we have:

$$\frac{\partial x^{(n)}}{\partial x^{(1)}} = \frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \dots \frac{\partial x^{(n)}}{\partial x^{(n-1)}} \quad (90)$$

□

Proposition 11. Consider a chain of: x a scalar, y a column vector, z a scalar. $z = z(y)$, $y_i = y_i(x)$, $i = 1, 2, \dots, n$. Apply the chain rule, we have

$$\frac{\partial z}{\partial x} = \frac{\partial y}{\partial x} \frac{\partial z}{\partial y} = \sum_{i=1}^n \frac{\partial y_i}{\partial x} \frac{\partial z}{\partial y_i} \quad (91)$$

Now we explain the intuition behind. x, z are both scalar, so we're basically calculating the derivative in ordinary calculus. Besides, we have a group of "bridge variables", y_i . $\frac{\partial y_i}{\partial x} \frac{\partial z}{\partial y_i}$ is just the result of applying scalar chain rule on the chain: $x \rightarrow y_i \rightarrow z$. The separate results of different bridge variables are additive! To see why I make this proposition, interested readers can refer to the comments in the corresponding L^AT_EX source file.

Example 8. Show the derivative of $(x - \mu)^T \Sigma^{-1} (x - \mu)$ to μ (for symmetric Σ^{-1}).

$$\frac{\partial [(x - \mu)^T \Sigma^{-1} (x - \mu)]}{\partial \mu} \quad (92)$$

$$= \frac{\partial [x - \mu]}{\partial \mu} \frac{\partial [(x - \mu)^T \Sigma^{-1} (x - \mu)]}{\partial [x - \mu]} \quad (93)$$

$$(example(4)) = \frac{\partial [x - \mu]}{\partial \mu} 2 \Sigma^{-1} (x - \mu) \quad (94)$$

$$(\mathrm{d}[x - \mu] = -I \mathrm{d}\mu) = -I 2 \Sigma^{-1} (x - \mu) \quad (95)$$

$$= -2 \Sigma^{-1} (x - \mu) \quad (96)$$

3 Application

3.1 The 2nd Induced Norm of Matrix

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

$$\|A\|_p = \max_x \frac{\|Ax\|_p}{\|x\|_p} \quad (97)$$

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:

$$\|A\|^2 = \max_x \frac{\|Ax\|^2}{\|x\|^2} \quad (98)$$

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

$$\|A\|^2 = \max_{x'} \frac{\|Ax'\|^2}{\|x'\|^2} = \max_x \frac{\|tAx\|^2}{\|tx\|^2} = \max_x \frac{t^2\|Ax\|^2}{t^2\|x\|^2} = \max_x \frac{\|Ax\|^2}{\|x\|^2} \quad (99)$$

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

$$\text{Maximize} \quad f(x) = \|Ax\|^2 \quad (100)$$

$$s.t. \quad \|x\|^2 = 1 \quad (101)$$

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

$$L(x) = f(x) - \lambda(\|x\|^2 - 1) \quad (102)$$

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

$$dL(x) = d\text{Tr}[L(x)] \quad (103)$$

$$= \text{Tr}[d(L(x))] \quad (104)$$

$$= \text{Tr}[d(x^T A^T A x - \lambda(x^T x - 1))] \quad (105)$$

$$= \text{Tr}[2x^T A^T A dx - \lambda(2x^T dx)] \quad (106)$$

$$= \text{Tr}[(2x^T A^T A - 2\lambda x^T)dx] \quad (107)$$

Next write out derivative:

$$\frac{\partial L}{\partial x} = 2A^T A x - 2\lambda x \quad (108)$$

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

$$(A^T A)x = \lambda x \quad (109)$$

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:

$$f(x) = x^T (A^T A x) = x^T (\lambda x) = \lambda \quad (110)$$

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

$$\|A\|^2 = \max_x f(x) = \lambda_{\max}(A^T A) \quad (111)$$

That is:

$$\|A\| = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max}(A) \quad (112)$$

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:

$$\lambda_{\max}^2(A) = \max_x \frac{\|Ax\|^2}{\|x\|^2} = \max_x \frac{x^T A^T A x}{x^T x} \quad (113)$$

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,

$$\frac{x^T A^T A x}{x^T x} \quad (114)$$

$$(v_k \text{ is an orthonormal set}) = \frac{\sum_i \lambda_i c_i^2}{\sum_i c_i^2} \quad (115)$$

$$\leq \frac{\sum_i \lambda_1 c_i^2}{\sum_i c_i^2} \quad (116)$$

$$= \lambda_1 \quad (117)$$

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

3.2 General Multivariate Gaussian Distribution

The first time I came across multivariate Gaussian distribution is in my sophomore year. However, at that time, the multivariate version is optional, and the text book only gives us the formula rather than explaining any intuition behind. I had trouble remembering the formula, since I don't know why it is the case.

During the Machine Learning course[12] I take recently, the rationale becomes clear to me. We'll start from the basic notion, generalize it to multivariate isotropic Gaussian, and then use matrix calculus to derive the most general version. The following content is adapted from the corresponding course notes and homework exercises.

We start from the univariate Gaussian[9]:

$$G(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (118)$$

The intuition behind is:

- Use one centroid to represent a set of samples coming from the same Gaussian mixture, which is denoted by μ .
- Allow the existence of error. We express the reverse notion of "error" by "likelihood". We want the density function describing the distribution that the farther away, the less likely a sample is drawn from the mixture. We want this likelihood to decrease in an accelerated manner. Negative exponential is one proper candidate, say $\exp\{-D(x, \mu)\}$, where D measures the distance between sample and centroid.
- To fit human's intuition, the "best" choice of distance measure is Euclidean distance, since we live in this Euclidean space. So far, we have $\exp\{-(x - \mu)^2\}$
- How fast the likelihood decreases should be controlled by a parameter, say $\exp\{\frac{-(x-\mu)^2}{\sigma}\}$. σ can also be thought to control the uncertainty.

Now we already get Gaussian distribution. The division by 2 in the exponent is only to simplify deductions. Writing σ^2 instead of σ is to align with some statistic quantities. The rest term is basically used to do normalization, which can be derived by extending the 1-D Gaussian integral to a 2-D area integral using Fubini's theorem[11]. Interested reader can refer to [10].

Now we're ready to extend the univariate version to multivariate version. The above four characteristics appear as simple interpretation to everyone who learned Gaussian distribution before. However, they're the real "axioms" behind. Now consider an isotropic Gaussian distribution and we start with perpendicular axes. The uncertainty along each axis is the same, so the distance measure is now $\|x - \mu\|_2^2 = (x - \mu)^T(x - \mu)$. The exponent is $\exp\{-\frac{1}{2\sigma^2}(x - \mu)^T(x - \mu)\}$. Integrating over the volume of x can be done by transforming it into iterated form using Fubini's theorem. Then we simply apply the method that we deal with univariate Gaussian integral. Now we have multivariate isotropic Gaussian distribution:

$$G(x|\mu, \sigma^2) = \frac{1}{(2\pi)^{d/2}\sigma^d} \exp\{-\frac{1}{2\sigma^2}(x - \mu)^T(x - \mu)\} \quad (119)$$

where d is the dimension of x, μ .

We are just one step away from a general Gaussian. Suppose we have a general Gaussian, whose covariance is not isotropic nor does it all along the standard Euclidean axes. Denote the sample from this Gaussian by ξ . We

can first apply a rotation on ξ to bring it back to the standard axes, which can be done by left multiplying an orthongonal matrix W^T . Then we scale each component of $W\xi$ by different ratio to make them isotropic, which can be done by left multiplying a diagonal matrix $\Lambda^{-0.5}$. The exponent -0.5 is simply to make later discussion convenient. Finally, we multiply it by σ to be able to control the uncertainty at each direction again. The transform is given by $x = \sigma\Lambda^{-0.5}W^T\xi$. Plugging this back to eqn(119), we derive the following exponent:

$$\exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^T(x - \mu)\right\} \quad (120)$$

$$= \exp\left\{-\frac{1}{2\sigma^2}(\sigma\Lambda^{-0.5}W^T\xi - \mu)^T(\sigma\Lambda^{-0.5}W^T\xi - \mu)\right\} \quad (121)$$

$$= \exp\left\{-\frac{1}{2\sigma^2}(\xi - \mu_\xi)^T\sigma^2W\Lambda^{-1}W^T(\xi - \mu_\xi)\right\} \quad (122)$$

$$= \exp\left\{-\frac{1}{2}(\xi - \mu_\xi)^T\Sigma^{-1}(\xi - \mu_\xi)\right\} \quad (123)$$

where we let $\mu_\xi = \mathbb{E}\xi$, and we plugged in the following result:

$$\mu = \mathbb{E}x = \mathbb{E}[\sigma\Lambda^{-0.5}W^T\xi] = \sigma\Lambda^{-0.5}W^T\mathbb{E}\xi = \sigma\Lambda^{-0.5}W^T\mu_\xi \quad (124)$$

Now we only need to find the normalization factor to make it a probability distribution. Note eqn(119) integrates to 1, which is:

$$\int G(x|\mu, \sigma^2)dx = \int \frac{1}{(2\pi)^{d/2}\sigma^d} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^T(x - \mu)\right\}dx = 1 \quad (125)$$

Transforming variable from x to ξ causes a scaling by absolute Jacobian, which we already calculated in example(7). That is:

$$\int \frac{1}{(2\pi)^{d/2}\sigma^d} \exp\left\{-\frac{1}{2}(\xi - \mu_\xi)^T\Sigma^{-1}(\xi - \mu_\xi)\right\}|J|d\xi = 1 \quad (126)$$

$$\int \frac{\sigma^d \det(\Sigma)^{-1/2}}{(2\pi)^{d/2}\sigma^d} \exp\left\{-\frac{1}{2}(\xi - \mu_\xi)^T\Sigma^{-1}(\xi - \mu_\xi)\right\}d\xi = 1 \quad (127)$$

$$\int \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left\{-\frac{1}{2}(\xi - \mu_\xi)^T\Sigma^{-1}(\xi - \mu_\xi)\right\}d\xi = 1 \quad (128)$$

The term inside integral is just the general Gaussian density we want to find:

$$G(x|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)^T\Sigma^{-1}(x - \mu)\right\} \quad (129)$$

3.3 Maximum Likelihood Estimation of Gaussian

Given N samples, $x_t, t = 1, 2, \dots, N$, independently indentially distributed that are drawn from the following Gaussian distribution:

$$G(x_t|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left\{-\frac{1}{2}(x_t - \mu)^T\Sigma^{-1}(x_t - \mu)\right\} \quad (130)$$

solve the paramters $\theta = \{\mu, \Sigma\}$, that maximize:

$$p(X|\theta) = \prod_{t=1}^N G(x_t|\theta) \quad (131)$$

It's more convenient to handle the log likelihood as defined below:

$$L = \ln p(X|\theta) = \sum_{t=1}^N \ln G(x_t|\theta) \quad (132)$$

We write each term of L out to facilitate further processing:

$$L = -\frac{Nd}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_t (x_t - \mu)^T \Sigma^{-1} (x_t - \mu) \quad (133)$$

Taking derivative of μ , the first two terms are gone and the third term is already handled by example(8):

$$\frac{\partial L}{\partial \mu} = \sum_t \Sigma^{-1} (x_t - \mu) \quad (134)$$

Let $\frac{\partial L}{\partial \mu} = 0$, we solve for $\mu = \frac{1}{N} \sum_t x_t$. Then take derivative of Σ . It eaiser to be handled using our trace schema:

$$dL = d[-\frac{N}{2} \ln |\Sigma|] + d[-\frac{1}{2} \sum_t (x_t - \mu)^T \Sigma^{-1} (x_t - \mu)] \quad (135)$$

The first term is:

$$d[-\frac{N}{2} \ln |\Sigma|] \quad (136)$$

$$= -\frac{N}{2} d \ln |\Sigma| \quad (137)$$

$$\text{(section(2.8))} = -\frac{N}{2} \text{Tr} [\Sigma^{-1} d\Sigma] \quad (138)$$

The second term is:

$$d[-\frac{1}{2} \sum_t (x_t - \mu)^T \Sigma^{-1} (x_t - \mu)] \quad (139)$$

$$= -\frac{1}{2} d \text{Tr} \left[\sum_t (x_t - \mu)^T \Sigma^{-1} (x_t - \mu) \right] \quad (140)$$

$$= -\frac{1}{2} d \text{Tr} \left[\sum_t (x_t - \mu)(x_t - \mu)^T \Sigma^{-1} \right] \quad (141)$$

$$\text{(example(5))} = -\frac{1}{2} \text{Tr} \left[\left[\sum_t (x_t - \mu)(x_t - \mu)^T \right] (-\Sigma^{-1} d\Sigma \Sigma^{-1}) \right] \quad (142)$$

$$= \frac{1}{2} \text{Tr} \left[\Sigma^{-1} \left[\sum_t (x_t - \mu)(x_t - \mu)^T \right] \Sigma^{-1} d\Sigma \right] \quad (143)$$

Then we have:

$$\frac{\partial L}{\partial \Sigma} = -\frac{N}{2}\Sigma^{-1} + \frac{1}{2}\Sigma^{-1}\left[\sum_t (x_t - \mu)(x_t - \mu)^T\right]\Sigma^{-1} \quad (144)$$

Let $\frac{\partial L}{\partial \Sigma} = 0$, we solve for $\Sigma = \frac{1}{N} \sum_t (x_t - \mu)(x_t - \mu)^T$.

4 Cheat Sheet

Acknowledgements

Thanks prof. XU, Lei's tutorial on matrix calculus[12]. Besides, the author also benefits 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>
- [7] Hessian Matrix, Wikipedia, http://en.wikipedia.org/wiki/Hessian_matrix
- [8] Jacobian Matrix, Wikipedia, http://en.wikipedia.org/wiki/Jacobian_matrix_and_determinant
- [9] Gaussian Distribution, Wikipedia, http://en.wikipedia.org/wiki/Normal_distribution
- [10] Gaussian Integral, Wikipedia, http://en.wikipedia.org/wiki/Gaussian_integral

- [11] Fubini's Theorem, Wikipedia, http://en.wikipedia.org/wiki/Fubini%27s_theorem
- [12] XU, Lei, 2012, Machine Learning Theory(CUHK CSCI5030).
- [13] Appendix from Introduction to Finite Element Methods book on University of Colorado at Boulder. <http://www.colorado.edu/engineering/cas/courses.d/IFEM.d/IFEM.AppD.d/IFEM.AppD.pdf>

Appendix