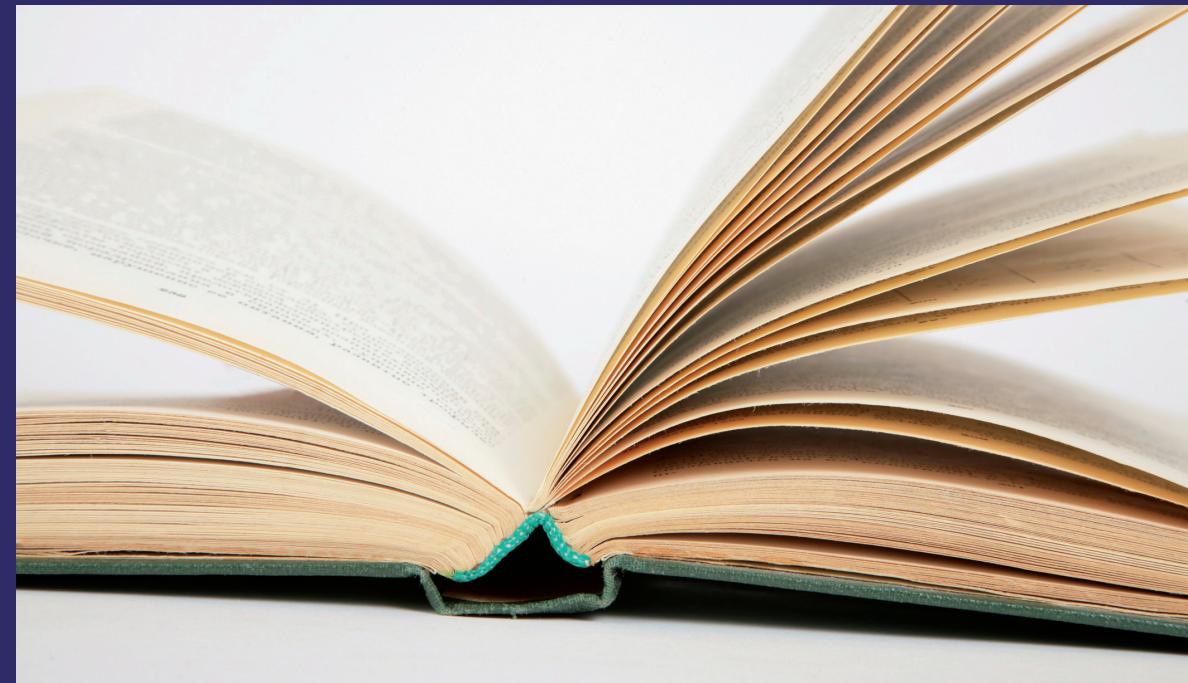


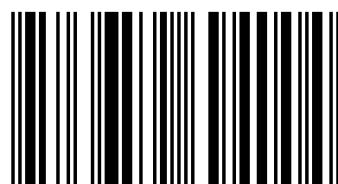
Statistics, multivariate data analysis and convex optimization are applied widely in many scientific domains and most analytical techniques are developed based on matrix analysis and matrix calculus because matrix is abstract representation of multivariate data. Although it is slightly confused for us to comprehend their concepts and theories, matrix analysis and calculus give us exciting results which enhance data analysis techniques to be more plentiful and accurate. So the report is survey of matrix analysis and calculus, which includes five main sections such as basic concepts, matrix analysis, matrix derivative, composite derivative, and applications of matrix. Matrix derivative and composite derivative are subjects of matrix calculus.



Loc Nguyen



Loc Nguyen is a director at Sunflower Soft Company, Vietnam. He is interested in computer science, statistics and mathematics. He serves as reviewer and editor in a wide range of international journals. Now he is a volunteer of Statistics Without Borders of American Statistics Association.



978-3-659-69400-4

Matrix Analysis and Calculus

Loc Nguyen

Matrix Analysis and Calculus

Loc Nguyen

Matrix Analysis and Calculus

LAP LAMBERT Academic Publishing

Impressum / Imprint

Bibliografische Information der Deutschen Nationalbibliothek: Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de> abrufbar.

Alle in diesem Buch genannten Marken und Produktnamen unterliegen warenzeichen-, marken- oder patentrechtlichem Schutz bzw. sind Warenzeichen oder eingetragene Warenzeichen der jeweiligen Inhaber. Die Wiedergabe von Marken, Produktnamen, Gebrauchsnamen, Handelsnamen, Warenbezeichnungen u.s.w. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutzgesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Bibliographic information published by the Deutsche Nationalbibliothek: The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at <http://dnb.d-nb.de>.

Any brand names and product names mentioned in this book are subject to trademark, brand or patent protection and are trademarks or registered trademarks of their respective holders. The use of brand names, product names, common names, trade names, product descriptions etc. even without a particular marking in this work is in no way to be construed to mean that such names may be regarded as unrestricted in respect of trademark and brand protection legislation and could thus be used by anyone.

Coverbild / Cover image: www.ingimage.com

Verlag / Publisher:

LAP LAMBERT Academic Publishing
ist ein Imprint der / is a trademark of
OmniScriptum GmbH & Co. KG
Heinrich-Böcking-Str. 6-8, 66121 Saarbrücken, Deutschland / Germany
Email: info@lap-publishing.com

Herstellung: siehe letzte Seite /

Printed at: see last page

ISBN: 978-3-659-69400-4

Copyright © 2015 OmniScriptum GmbH & Co. KG

Alle Rechte vorbehalten. / All rights reserved. Saarbrücken 2015

Matrix Analysis and Calculus

Loc Nguyen

Vietnam Institute of Mathematics

Abstract

Statistics, multivariate data analysis and convex optimization are applied widely in many scientific domains and most analytical techniques are developed based on matrix analysis and matrix calculus because matrix is abstract representation of multivariate data. Although it is slightly confused for us to comprehend their concepts and theories, matrix analysis and calculus give us exciting results which enhance data analysis techniques to be more plentiful and accurate. So the report is survey of matrix analysis and calculus, which includes five main sections such as basic concepts, matrix analysis, matrix derivative, composite derivative, and applications of matrix. Matrix derivative and composite derivative are subjects of matrix calculus.

Contents

1. Basic concepts	1
2. Matrix analysis.....	10
3. Matrix derivative.....	21
4. Composite derivative	32
5. Some applications of matrix analysis and calculus.....	50
6. Conclusion	59
Acknowledgement	60
Bibliography.....	60

1. Basic concepts

We begin matrix analysis and matrix calculus with some basic concepts of matrix. Before discussing main subjects, there are some conventions. Firstly, if there is no additional explanation, normal letters denote scalar variables and numbers; vectors are denoted by bold letters and matrices are denoted as bold and uppercase letters. For example, letters x , \mathbf{x} and X indicate a scalar variable x , a vector \mathbf{x} and a matrix X . Note that term “scalar” refers real or complex number in opposite to vector and matrix. All constants are denoted by lowercase letters, for example, letters c , a and A express a scalar constant, a vector constant and a matrix constant, respectively and so their difference is based on study context. If there is the mix of constant vectors and constant scalars in the same expression, Greek lowercase letters such as α , β and γ often denote constant vectors and Latin lowercase letters such as a , b and c often denote scalar constants. Analysis spaces or algebra structures are denoted by letter-like symbols such as complex number field \mathbb{C} , real number field \mathbb{R} , rational number field \mathbb{Q} , integer ring \mathbb{Z} , set of natural number \mathbb{N} , vector space \mathbb{V} , n -dimension vector space over real number field \mathbb{R}^n . The default vector space is \mathbb{R}^n .

Vector and matrix properties

Vector \mathbf{x} is a range of n values denoted $\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ or $(x_1, x_2, \dots, x_n)^T$ called n -dimension vector where T

denotes transpose operation which changes row to column and otherwise. Please pay attention to transpose operation because it is used over the whole of report. There are two kinds of vector such

as aforementioned column vector $\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ and row vector (x_1, x_2, \dots, x_n) . By default, \mathbf{x} denotes

column vector and \mathbf{x}^T denotes row vector. Each x_i is partial value (partial component, partial element, or partial scalar) of vector \mathbf{x} . If \mathbf{x} is random variable, x_i is also random variable. Simple vector operations include transposition, addition, subtraction, scalar multiplication and scalar division.

$$\begin{aligned} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}^T &= (x_1, x_2, \dots, x_n) \text{ and } (x_1, x_2, \dots, x_n)^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \\ \mathbf{x} \pm \mathbf{y} &= \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \pm \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 \pm y_1 \\ x_2 \pm y_2 \\ \vdots \\ x_n \pm y_n \end{pmatrix} \\ c\mathbf{x} &= c \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} cx_1 \\ cx_2 \\ \vdots \\ cx_n \end{pmatrix} \text{ and } \frac{\mathbf{x}}{c} = \frac{1}{c} \mathbf{x} = \begin{pmatrix} x_1/c \\ x_2/c \\ \vdots \\ x_n/c \end{pmatrix} \end{aligned}$$

The *dot product* of two vectors \mathbf{x} and \mathbf{y} denoted $\mathbf{x} \cdot \mathbf{y}$ or $\mathbf{x}\mathbf{y}^T$ is a scalar value which is sum of multiplications of their components. The dot product is also called inner product or scalar product.

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i$$

The length also called norm, module, or magnitude of vector \mathbf{x} is:

$$|\mathbf{x}| = \|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}$$

Both notations $|.|$ and $\|.\|$ are used to indicate length of vector but the notation $|.|$ is used more frequently. We also use the notation $|.|$ to indicate absolute value of scalar and determinant of matrix. Normalized vector is defined as vector whose length is 1. Arbitrary vector is normalized by dividing itself by its length.

$$\text{normalize}(\mathbf{x}) = \frac{\mathbf{x}}{|\mathbf{x}|} = \begin{pmatrix} x_1/|\mathbf{x}| \\ x_2/|\mathbf{x}| \\ \vdots \\ x_n/|\mathbf{x}| \end{pmatrix}$$

The cosine of angle between two vectors \mathbf{x} and \mathbf{y} is the ratio of its product to multiplication of their lengths.

$$\cos(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$$

Cosine ranges in interval $[-1, 1]$. The *cross product* or outer product of two vectors \mathbf{x} and \mathbf{y} denoted $\mathbf{x} \times \mathbf{y}$ is the vector that is perpendicular to both \mathbf{x} and \mathbf{y} according to right-hand rule. The length of cross product vector is:

$$|\mathbf{x} \times \mathbf{y}| = |\mathbf{x}| |\mathbf{y}| \sin(\mathbf{x}, \mathbf{y})$$

Where $\sin(\mathbf{x}, \mathbf{y})$ is the sine of angle between two vectors \mathbf{x} and \mathbf{y} , which ranges in $[0, \pi]$ according to right-hand rule, hence $\sin(\mathbf{x}, \mathbf{y})$ is always greater than or equal to 0 and is calculated by formula $\sin(\mathbf{x}, \mathbf{y}) = \sqrt{1 - \cos^2(\mathbf{x}, \mathbf{y})} = \sqrt{1 - \left(\frac{\mathbf{x}^T \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|} \right)^2}$. In 3-dimensional vector space, the cross product can be computed via determinant of second order square matrix. Given vectors $\mathbf{x} = (x_1, x_2, x_3)$ and $\mathbf{y} = (y_1, y_2, y_3)$, their cross product $\mathbf{z} = \mathbf{x} \times \mathbf{y}$ is determined as follows:

$$\mathbf{z} = \mathbf{x} \times \mathbf{y} = \left(\begin{vmatrix} x_2 & x_3 \\ y_2 & y_3 \end{vmatrix}, \begin{vmatrix} x_3 & x_1 \\ y_3 & y_1 \end{vmatrix}, \begin{vmatrix} x_1 & x_2 \\ y_1 & y_2 \end{vmatrix} \right)$$

Where $\begin{vmatrix} x_i & x_j \\ y_i & y_j \end{vmatrix} = (x_i y_j - x_j y_i)$ is determinant of square matrix $\begin{pmatrix} x_i & x_j \\ y_i & y_j \end{pmatrix}$, which is discussed later.

Note that the semantic of cross product $\mathbf{x} \times \mathbf{y} = \mathbf{z}$ is the same to the semantic of normal multiplication like $2 * 3 = 6$.

According to (Baker K. , 2013), given an integer $k > 1$, the k^{th} power of vector \mathbf{x} is defined as below:

$$\begin{cases} \mathbf{x}^k = \prod_{i=1}^{k/2} \mathbf{x}^T \mathbf{x} & \text{if } k \text{ is even} \\ \mathbf{x}^k = \mathbf{x} \times \prod_{i=1}^{(k-1)/2} \mathbf{x}^T \mathbf{x} & \text{if } k \text{ is odd} \end{cases}$$

Where Π denotes cross product or scalar product. So we infer that \mathbf{x}^k is scalar if k is even and \mathbf{x}^k is vector if k is odd. We have summarization table for $k = 2, 3, 4, 5$ as follows (Baker M. J., n.d.):

Power	Value	Type
\mathbf{x}^2	$\mathbf{x} \cdot \mathbf{x}$	Scalar
\mathbf{x}^3	$(\mathbf{x} \cdot \mathbf{x}) \times \mathbf{x}$	Vector
\mathbf{x}^4	$(\mathbf{x} \cdot \mathbf{x}) \times (\mathbf{x} \cdot \mathbf{x})$	Scalar
\mathbf{x}^5	$(\mathbf{x} \cdot \mathbf{x}) \times (\mathbf{x} \cdot \mathbf{x}) \times \mathbf{x}$	Vector

Note, given vector \mathbf{x} and scalars a, b , we have:

$$a \times \mathbf{x} = a\mathbf{x} \text{ and } a \times b = ab$$

Matrix \mathbf{A} is a table including m rows and n columns, whose cells or elements are scalar values.

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}$$

Matrix \mathbf{A} denoted (a_{ij}) , $\mathbf{A}(m,n)$, $\mathbf{A}_{m,n}$, or $m \times n$ matrix \mathbf{A} can be considered as a set of m row vectors or a set of n column vectors.

$$\begin{aligned} \mathbf{A}(m,n) &= ((a_{11}, a_{12}, \dots, a_{1n}), (a_{21}, a_{22}, \dots, a_{2n}), \dots, (a_{m1}, a_{m2}, \dots, a_{mn})) \\ &= \left(\begin{pmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{m1} \end{pmatrix}, \begin{pmatrix} a_{12} \\ a_{22} \\ \vdots \\ a_{m2} \end{pmatrix}, \dots, \begin{pmatrix} a_{1n} \\ a_{2n} \\ \vdots \\ a_{mn} \end{pmatrix} \right) \end{aligned}$$

Matrix $\mathbf{A}(n,n)$ having the same number of rows and columns is called square matrix. Square matrix is very popular in computational application. Matrix $\mathbf{A}(m,n)$ having different number of rows and columns is often called rectangle or non-square matrix. Vector can be considered as 1-row or 1-column matrix, so $\mathbf{X}(1,n)$ denote row vector and $\mathbf{X}(n,1)$ denotes column vector. By default, vector \mathbf{X} is column vector if there is no additional note. Given matrices $\mathbf{A}(m,n)$ and $\mathbf{B}(n,p)$ and scalar constant c , following are matrix operations including transposition, multiplication, addition, subtraction and scalar multiplication.

$\mathbf{A}^T = (a_{ji})$ where T is transpose operation which changes row to column and otherwise.

$$\mathbf{A} + \mathbf{B} = (a_{ij}) + (b_{ij})$$

$$\mathbf{A} - \mathbf{B} = (a_{ij}) - (b_{ij})$$

$$c\mathbf{A} = (ca_{ij})$$

$$\mathbf{AB} = \mathbf{C}(m,p) = (c_{ij}) = (\sum_{k=1}^n a_{ik} b_{jk}) \text{ where the number of columns of } \mathbf{A} \text{ is equal to the number of}$$

rows of \mathbf{B} . This is *multiplication condition*.

Given $n \times n$ square matrix \mathbf{A} , the n^{th} power of \mathbf{A} is defined as $\mathbf{A}^k = \prod_{i=1}^k \mathbf{A}$. It is convention that $\mathbf{A}^0 = \mathbf{I}_n$ where $\mathbf{I}_n = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}$ is identity matrix whose diagonal elements are 1 and remaining elements are 0. Following are properties of matrix operations.

$$\begin{aligned}(A^T)^T &= A \\ (AB)^T &= B^T A^T \\ A + B &= B + A \\ A(B + C) &= AB + AC \\ A(BC) &= (AB)C\end{aligned}$$

Note that matrix multiplication is not commutative. Following is the list of some special vectors and matrices (Härdle & Simar, 2013, p. 59).

Name	Definition	Notation	Example
Scalar	Scalar value is considered as 1-element vector or matrix	c	1
Column vector	$(a_1, a_2, \dots, a_n)^T$	a	$\begin{pmatrix} 1 \\ 2 \end{pmatrix}$
Row vector	(a_1, a_2, \dots, a_n)	a^T	$(1, 2)$
Vector of ones	$(1, 1, \dots, 1)^T$	$\mathbf{1}_n$	$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$
Vector of zeros or zero vector	$(0, 0, \dots, 0)^T$	$\mathbf{0}_n$	$\begin{pmatrix} 0 \\ 0 \end{pmatrix}$
Square matrix	Having the same number of rows and columns	$A(n \times n)$	$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$
Diagonal matrix	Square matrix and $a_{ij} = 0, \forall i \neq j$	$\text{diag}(a_{ii})$	$\begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$
Identity matrix	Diagonal matrix and $a_{ij} = 0, \forall i \neq j$ and $a_{ii} = 1$	\mathbf{I}_n	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
Unit matrix	Square matrix and $a_{ij} = 1$	$\mathbf{1}_n \mathbf{1}_n^T$	$\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$
Null matrix or zero matrix	Square matrix and $a_{ij} = 0$	$\mathbf{0}_n \mathbf{0}_n^T$	$\begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$
Symmetric matrix	Square matrix and $a_{ij} = a_{ji}$		$\begin{pmatrix} 1 & 3 \\ 3 & 2 \end{pmatrix}$
Upper triangular matrix	Square matrix and $a_{ij} = 0, \forall i < j$		$\begin{pmatrix} 1 & 4 & 5 \\ 0 & 2 & 6 \\ 0 & 0 & 3 \end{pmatrix}$
Idempotent matrix	Square matrix and $\mathbf{AA} = \mathbf{A}$		$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1/2 & 1/2 \\ 0 & 1/2 & 1/2 \end{pmatrix}$
Nilpotent matrix	Square matrix and $\underbrace{\mathbf{AA} \dots \mathbf{A}}_{k \text{ times}} = \mathbf{0}_n \mathbf{0}_n^T$		$\begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$

Orthogonal matrix	Square matrix and $\mathbf{A}\mathbf{A}^T = \mathbf{I}_n = \mathbf{A}^T\mathbf{A}$		$\begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$
-------------------	---	--	---

The notation $\mathbf{0}$ can indicate both zero vector $\mathbf{0}_n$ and zero matrix $\mathbf{0}_n\mathbf{0}_n^T$ according to study context. The notation \mathbf{I} without n -index implicates identity $n \times n$ matrix. Moreover diagonal matrix – an important kind of square matrix with regard to matrix analysis can be denoted as a composition of column vectors such as $\mathbf{A}(n,n) = (a_{ij}) = diag(\lambda_1, \lambda_2, \dots, \lambda_n)$ where λ_i are n -dimensional vector whose elements are zero except i^{th} element. Following is an example of diagonal matrix.

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix} = diag(\lambda_1, \lambda_2, \lambda_3) \text{ where } \lambda_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \lambda_2 = \begin{pmatrix} 0 \\ 2 \\ 0 \end{pmatrix}, \lambda_3 = \begin{pmatrix} 0 \\ 0 \\ 3 \end{pmatrix}$$

We assume that all vectors and matrices have n dimensions if there is no additional explanation. Now we survey characteristics of matrix. The rank of $m \times n$ matrix \mathbf{A} denoted $rank(\mathbf{A})$ is the maximum number of linear independent rows or columns, $rank(\mathbf{A}) \leq \min(m, n)$. Non-singular or invertible matrix is square matrix and all its rows (columns) are independent; so the rank of non-singular is its number of rows (columns).

The trace of square matrix \mathbf{A} denoted $tr(\mathbf{A})$ is the sum of all its diagonal, $tr(\mathbf{A}) = \sum a_{ii}$. Following are properties of *trace* and *rank* (Härdle & Simar, 2013, p. 62).

$$tr(\mathbf{A} + \mathbf{B}) = tr(\mathbf{A}) + tr(\mathbf{B})$$

$$tr(c\mathbf{A}) = c tr(\mathbf{A})$$

$$tr(\mathbf{AB}) = tr(\mathbf{BA}) \text{ given } \mathbf{A}(m,n) \text{ and } \mathbf{B}(n,m)$$

$$tr(\mathbf{ABC}) = tr(\mathbf{BCA}) = tr(\mathbf{CAB}) \text{ given } \mathbf{A}(m,n), \mathbf{B}(n,p) \text{ and } \mathbf{C}(p,m)$$

$$tr(\mathbf{ABCD} \dots \mathbf{Z} \dots) = tr(\dots \mathbf{Z} \dots \mathbf{DCBA}) \text{ when multiplication condition is satisfied.}$$

$$tr(\mathbf{A}) = tr(\mathbf{A}^T)$$

$$tr(\mathbf{A}^T \mathbf{B}) = tr(\mathbf{AB}^T)$$

$$\mathbf{a}^T \mathbf{a} = tr(\mathbf{aa}^T) \text{ where } \mathbf{a} \text{ is a vector}$$

$$rank(\mathbf{A}) \leq \min(m, n) \text{ given } \mathbf{A}(m,n)$$

$$rank(\mathbf{A}) \geq 0$$

$$rank(\mathbf{A}) = rank(\mathbf{A}^T)$$

$$rank(\mathbf{A}(n,n)) = n \text{ if } \mathbf{A} \text{ is non-singular}$$

$$rank(\mathbf{A}^T \mathbf{A}) = rank(\mathbf{AA}^T) = rank(\mathbf{A})$$

$$rank(\mathbf{A} + \mathbf{B}) = rank(\mathbf{A}) + rank(\mathbf{B})$$

$$rank(\mathbf{AB}) \leq \min(rank(\mathbf{A}), rank(\mathbf{B}))$$

$$rank(\mathbf{ABC}) = rank(\mathbf{B}) \text{ given } \mathbf{A} \text{ and } \mathbf{C} \text{ are non-singular.}$$

Please pay attention to trace operator because it is important linear operator which relates to matrix derivative.

Matrix determinant and inverse

Given a square matrix $\mathbf{A}(n,n)$, determinant of matrix \mathbf{A} denoted $det(\mathbf{A})$ or $|\mathbf{A}|$ is defined as the sum over all permutations σ (s) of indexes $\{1, 2, \dots, n\}$. The i^{th} value of permutation σ is denoted

permutations σ_i . The set of all permutations σ (s) given n indexes $\{1, 2, \dots, n\}$ is denoted S_n . S_n has $n!$ elements (n factorial elements).

$$|\mathbf{A}| = \det(\mathbf{A}) = \sum_{\sigma \in S_n} sgn(\sigma) \prod_{i=1}^n a_{i\sigma_i}$$

Formula above is Leibniz formula where $sgn(\sigma)$ denotes the sign of permutation σ .

$$\begin{cases} sgn(\sigma) = 1 & \text{if } \sigma \text{ is even} \\ sgn(\sigma) = -1 & \text{if } \sigma \text{ is odd} \end{cases}$$

Permutation σ is even (or odd) when the new permutation can be obtained by the even (or odd, respectively) number of switches of numbers. For example, given indexes $\{1, 2, 3\}$, the permutation (123) is even because there is 0 number of switches and the permutation (132) is even because there is 1 number of switches. According to Leibniz formula, the determinants of 2×2 and 3×3 matrices are:

$$|\mathbf{A}(2,2)| = a_{11}a_{22} - a_{12}a_{21}$$

$$|\mathbf{A}(3,3)| = a_{11}a_{22}a_{33} - a_{11}a_{23}a_{32} - a_{12}a_{21}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32} - a_{13}a_{22}a_{31}$$

Leibniz formula gets a huge of operators due to $n!$ permutations. So matrix determinant is computed more effectively by applying Laplace expansion. Given element a_{ij} of square matrix $\mathbf{A}(n,n)$, the algebra complement denoted \mathbf{M}_{ij} is the sub-matrix including $n - 1$ rows and $n - 1$ columns, which is created by removing i^{th} row and j^{th} column from \mathbf{A} (Nguyen, 1999, p. 130).

$$\mathbf{M}_{ij} = \begin{pmatrix} a_{11} & \cdots & a_{1(j-1)} & a_{1(j+1)} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_{(i-1)1} & \cdots & a_{(i-1)(j-1)} & a_{(i-1)(j+1)} & \cdots & a_{(i-1)n} \\ a_{(i+1)1} & \cdots & a_{(i+1)(j-1)} & a_{(i+1)(j+1)} & \cdots & a_{(i+1)n} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{n(j-1)} & a_{n(j+1)} & \cdots & a_{nn} \end{pmatrix}$$

\mathbf{M}_{ij} is called algebra complement or complement matrix of \mathbf{A} at entry (i, j) . By the similar way, the elementary matrix denoted \mathbf{E}_{ij} of \mathbf{A} at entry (i, j) is defined as the matrix having all zero elements except element a_{ij} .

$$\mathbf{E}_{ij} = \begin{pmatrix} 0 & 0 & \cdots & 0 \\ 0 & 0 & \vdots & 0 \\ \vdots & \vdots & a_{ij} & \vdots \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

According to Laplace extension, determinant of $n \times n$ square matrix \mathbf{A} is computed according to determinants of algebra complements given arbitrary i^{th} row.

$$|\mathbf{A}| = \sum_{j=1}^n (-1)^{i+j} a_{ij} |\mathbf{M}_{ij}|$$

Formula above is recursive formula with regard to column expansion. According to row expansion, determinant of square matrix $\mathbf{A}(n,n)$ given arbitrary j^{th} column is:

$$|\mathbf{A}| = \sum_{i=1}^n (-1)^{i+j} a_{ij} |\mathbf{M}_{ij}|$$

Note that the expression $(-1)^{i+j}/|\mathbf{M}_{ij}|$ is called cofactor c_{ij} (Härdle & Simar, 2013, p. 60). Cofactor matrix of \mathbf{A} denoted $\text{cofactor}(\mathbf{A})$ is defined as matrix whose elements are cofactor c_{ij} (s). The adjoint matrix of \mathbf{A} denoted $\text{adj}(\mathbf{A})$ is the transposition of cofactor matrix of \mathbf{A} .

$$\text{adj}(\mathbf{A}) = (\text{cofactor}(\mathbf{A}))^T = (c_{ji})$$

Given square matrix $\mathbf{A}(n,n)$ and $|\mathbf{A}| \neq 0$, the inverse of \mathbf{A} denoted \mathbf{A}^{-1} is the one that multiplication of \mathbf{A} and itself is equal to identical matrix.

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_n$$

The inverse \mathbf{A}^{-1} is determined according to adjoint matrix.

$$\mathbf{A}^{-1} = \frac{1}{|\mathbf{A}|} \text{adj}(\mathbf{A})$$

If there is existence of the inverse of \mathbf{A} , in other words $|\mathbf{A}| \neq 0$, then \mathbf{A} is invertible or non-singular. If \mathbf{A} is orthogonal matrix, it is very easy to determine the inverse because of $\mathbf{A}^{-1} = \mathbf{A}^T$. The generalized inverse (G-inverse) of matrix \mathbf{A} denoted \mathbf{A}^- is the one satisfying following condition (Härdle & Simar, 2013, p. 60).

$$\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}$$

Note that the inverse \mathbf{A}^{-1} is a concrete case of G-inverse \mathbf{A}^- and \mathbf{A}^- is the general concept of the inverse \mathbf{A}^{-1} , hence, if \mathbf{A}^{-1} exists then $\mathbf{A}^- = \mathbf{A}^{-1}$. The G-inverse is also called pseudo unique and always exists even though \mathbf{A}^{-1} does not exist or \mathbf{A} is not square matrix.

A typical application of matrix determinant and inverse is to solve set of linear equations. Matrix algebra is originated from research how to solve a set of linear equations but after that, theories of abstract algebra such as group, ring, field, module and vector space are applied into researching it and so matrix algebra is generalized as a large domain with a lot of applications. Given a set of n linear equations whose a_{ij} are coefficients and x_{ij} are unknowns need solved (Nguyen, 1999, pp. 136-138).

$$\begin{cases} a_{11}x_1 + a_{12}x_1 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_1 + \cdots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_1 + \cdots + a_{nn}x_n = b_n \end{cases}$$

The set of equations above can be solved by Gaussian method including primary transformations such as changing order of two equations, multiplying a equation by scalar and adding to a equation by a linear combination of other equations so as to transform the set of equations into the form that is as simple as possible in order to find out unknowns.

In the report, we use another method called Cramer's rule to solve the set of linear equations. The set of n linear equations is re-written in matrix form.

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

Where $\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \vdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{nn} & a_{nn} \end{pmatrix}$ is $n \times n$ matrix and

$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ is unknown vector and

$\mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$ is free – coefficient vector

We have a convention that $\mathbf{A}\mathbf{x} = \mathbf{b}$ is vector equation and \mathbf{x} is unknown vector and $\mathbf{x}_0 = (x_{01}, x_{02}, \dots, x_{0n})$ is the solution vector and x_{0n} is called partial solution. Suppose matrix \mathbf{A} is non-singular and so its inverse \mathbf{A}^{-1} exists, Cramer's rule states that the number of partial solutions is the same to the number of partial equations; it means that matrix equation has unique solution \mathbf{x}_0 which has n partial solutions x_{0j} (s) and each solution x_{0j} is found out via determinants (Nguyen, 1999, p. 137):

$$x_{0j} = \frac{|A_j|}{|A|}$$

Where A_j is the matrix constructed by replacing j^{th} column in matrix by free-coefficient vector \mathbf{b} . For example, it is required to solve following equations (Nguyen, 1999, p. 138):

$$\begin{cases} x + y + 3z + 4t = -3 \\ x + y + 5z + 2t = 1 \\ 2x + y + 3z + 2t = -3 \\ 2x + 3y + 11z + 5t = 2 \end{cases}$$

The set of equations is written in matrix form $\mathbf{A}\mathbf{x} = \mathbf{b}$:

$$\begin{pmatrix} 1 & 1 & 3 & 4 \\ 1 & 1 & 5 & 2 \\ 2 & 1 & 3 & 2 \\ 2 & 3 & 11 & 5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ t \end{pmatrix} = \begin{pmatrix} -3 \\ 1 \\ -3 \\ 2 \end{pmatrix}$$

Determinant of \mathbf{A} is:

$$\begin{vmatrix} 1 & 1 & 3 & 4 \\ 1 & 1 & 5 & 2 \\ 2 & 1 & 3 & 2 \\ 2 & 3 & 11 & 5 \end{vmatrix} = -14$$

Solutions of these equations are totally solved as follows:

$$x = \frac{\begin{vmatrix} -3 & 1 & 3 & 4 \\ 1 & 1 & 5 & 2 \\ -3 & 1 & 3 & 2 \\ 2 & 3 & 11 & 5 \end{vmatrix}}{-14} = -2, y = \frac{\begin{vmatrix} 1 & -3 & 3 & 4 \\ 1 & 1 & 5 & 2 \\ 2 & -3 & 3 & 2 \\ 2 & 2 & 11 & 5 \end{vmatrix}}{-14} = 0,$$

$$z = \frac{\begin{vmatrix} 1 & 1 & -3 & 4 \\ 1 & 1 & 1 & 2 \\ 2 & 1 & -3 & 2 \\ 2 & 3 & 2 & 5 \end{vmatrix}}{-14} = 1, y = \frac{\begin{vmatrix} 1 & 1 & 3 & -3 \\ 1 & 1 & 5 & 1 \\ 2 & 1 & 3 & -3 \\ 2 & 3 & 11 & 2 \end{vmatrix}}{-14} = -1$$

Following are properties of matrix determinant and matrix inverse with note that determinant exists if and only if matrix is invertible.

$$|\mathbf{I}_n| = 1$$

$$|\mathbf{A}^T| = |\mathbf{A}|$$

$$|\mathbf{A}^{-1}| = |\mathbf{A}|^{-1}$$

$$|\mathbf{A}| \neq 0 \Leftrightarrow \text{rank}(\mathbf{A}) = n \Leftrightarrow \mathbf{A} \text{ non-singular} \Leftrightarrow \mathbf{A} \text{ invertible}$$

$$|\mathbf{AB}| = |\mathbf{A}| |\mathbf{B}|$$

$$|c\mathbf{A}| = c^n |\mathbf{A}| \text{ where } c \text{ is scalar constant.}$$

$$|\mathbf{A}^n| = |\mathbf{A}|^n$$

$$|\mathbf{A}| = \prod_{i=1}^n a_{ii} \text{ if } \mathbf{A} \text{ is triangular matrix.}$$

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$$

$$(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T$$

$$|\mathbf{I} + \mathbf{ab}^T| = 1 + \mathbf{a}^T \mathbf{b} \text{ where } \mathbf{I} \text{ is } n \times n \text{ identity matrix and } \mathbf{a}, \mathbf{b} \text{ are vectors.}$$

2. Matrix analysis

The first section introduces basic concepts relevant to matrix algebra. Pure matrix algebra uses theories of abstract algebra such as group, ring, module, and field to explain and solve theoretical problems of matrix. Matrix analysis and matrix calculus mentioned in this report are typical applications of matrix algebra. Matrix analysis focuses on processing and analyzing multivariate data while matrix calculus focuses on derivative and differential with regard to matrix. Matrix analysis is very necessary to matrix calculus, in which spectrum decomposition is the most important. This section includes four parts such as spectrum decomposition, singular decomposition, quadratic form and partitioned matrix. Its main contents are extracted from the book “Applied Multivariate Statistical Analysis” by authors (Härdle & Simar, 2013) which is a valued document that readers should read.

Spectrum decomposition

Matrix analysis consists of decomposing matrix techniques which transform a complex matrix into a set of simple matrices because it is easy to process partially on simple matrices instead of whole complex matrix. This is problem of spectrum decomposition. Matrix analysis starts with concepts of eigenvalues and eigenvectors. Given square matrix $A(n,n)$, if there is a scalar value $\lambda > 0$ and a vector u such that

$$A = \lambda u$$

Then λ and u are eigenvalue and eigenvector of matrix A , respectively. If $|A| \neq 0$ then there are n eigenvalues and n respective eigenvectors. At that time, the vector space determined by matrix A is decomposed into n disjoint sub-spaces and each sub-space is specified by a pair of eigenvalue and eigenvector. There is a question how to find out n eigenvalues and n respective eigenvectors and so, eigenvalues are solutions of following equation:

$$|A - \lambda I_n| = 0 \text{ where } I_n \text{ is } n,n \text{ identity matrix.}$$

Determinant $|A - \lambda I_n|$ is expanded as a n^{th} order polynomial which has n solutions $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ with attention that λ_i (s) should be in descending ordering, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. Such polynomial is called characteristic polynomial of matrix A with respect to (w.r.t) variable scalar λ , denoted $P(\lambda) = |A - \lambda I_n|$. Suppose a solution of $P(\lambda) = |A - \lambda I_n|$ is λ_i , the respective eigenvector u_i is vector solution of following set of linear equations:

$$(A - \lambda_i I_n)x = 0 \Leftrightarrow \begin{pmatrix} a_{11} - \lambda_{i1} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda_{i2} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \lambda_{in} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$
$$\Leftrightarrow \begin{cases} (a_{11} - \lambda_{i1})x_1 + a_{12}x_2 + \dots + a_{1n}x_n = 0 \\ a_{21}x_1 + (a_{22} - \lambda_{i2})x_2 + \dots + a_{2n}x_n = 0 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + (a_{nn} - \lambda_{in})x_n = 0 \end{cases}$$

In solution space of equation $(\mathbf{A} - \lambda \mathbf{I}_n)\mathbf{x} = \mathbf{0}$, we should choose eigenvectors \mathbf{u}_i so that they are mutually orthogonal and all of eigenvectors \mathbf{u}_i must be normalized. Maybe some eigenvectors in solution space are not mutually orthogonal but we assume that the mentioned eigenvectors are mutually orthogonal in whole report. Let $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ be the eigenvalue matrix and let \mathbf{U} be the orthogonal matrix created by n eigenvectors. \mathbf{U} is called eigenvector matrix. Both Λ and \mathbf{U} are spectrums of \mathbf{A} .

$$\Lambda = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{pmatrix}$$

$$\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n) = \begin{pmatrix} u_{11} & u_{21} & \dots & u_{n1} \\ u_{12} & u_{22} & \dots & u_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ u_{1n} & u_{2n} & \dots & u_{nn} \end{pmatrix}$$

where $\mathbf{u}_i = \begin{pmatrix} u_{i1} \\ u_{i2} \\ \vdots \\ u_{in} \end{pmatrix}$ is a column normalized eigenvector

The *Jordan decomposition theorem* (Härdle & Simar, 2013, p. 63) states that

$$\mathbf{A} = \mathbf{U}\Lambda\mathbf{U}^{-1} = \mathbf{U}\Lambda\mathbf{U}^T = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

This theorem is the most important theorem in matrix analysis and it is a base of many decomposition techniques. Jordan decomposition exists if and only if \mathbf{A} is square matrix and non-singular, $|\mathbf{A}| \neq 0$ and such matrix \mathbf{A} is call *diagonalizable* or diagonalized matrix. In general, Jordan decomposition is essentially spectrum decomposition. Author (Hoang, 2012) gives an example for illustrating Jordan decomposition or diagonalizing matrix as follows:

Suppose $\mathbf{A} = \begin{pmatrix} 3 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & 5 \end{pmatrix}$ and its characteristic polynomial $P(\lambda)$ is:

$$P(\lambda) = |\mathbf{A} - \lambda \mathbf{I}_3| = \begin{vmatrix} 3 - \lambda & -2 & 0 \\ -2 & 3 - \lambda & 0 \\ 0 & 0 & 5 - \lambda \end{vmatrix} = (1 - \lambda)(5 - \lambda)^2$$

Solving equation $P(\lambda) = 0$ results out three solutions λ_1, λ_2 and λ_3 which are eigenvalues.

$$P(\lambda) = 0 \Leftrightarrow \begin{cases} \lambda_1 = 1 \\ \lambda_2 = \lambda_3 = 5 \end{cases} \text{ (dual)}$$

Substituting λ_1, λ_2 and λ_3 into the equation $(\mathbf{A} - \lambda \mathbf{I}_3)\mathbf{x} = \mathbf{0}$, we get eigenvectors $\mathbf{u}_1, \mathbf{u}_2$ and \mathbf{u}_3 as solutions and normalize them right after.

$$(\mathbf{A} - \lambda_1 \mathbf{I}_3) \mathbf{x} = \mathbf{0} \Leftrightarrow \left(\begin{pmatrix} 3 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & 5 \end{pmatrix} - 1 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = x_1 \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \Rightarrow \mathbf{u}_1 = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \\ 0 \end{pmatrix}$$

$$(\mathbf{A} - \lambda_2 \mathbf{I}_3) \mathbf{x} = \mathbf{0} \Leftrightarrow \left(\begin{pmatrix} 3 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & 5 \end{pmatrix} - 5 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = x_1 \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} + x_3 \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \Rightarrow \begin{cases} \mathbf{u}_2 = \begin{pmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \\ 0 \end{pmatrix} \\ \mathbf{u}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \end{cases}$$

In general, we have eigenvalue matrix Λ and eigenvector matrix \mathbf{U} as following.

$$\mathbf{U} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 0 & 0 & 1 \end{pmatrix} \text{ and } \Lambda = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{pmatrix}$$

It is easy for us to validate spectrum decomposition:

$$\begin{aligned} \begin{pmatrix} 3 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & 5 \end{pmatrix} &= \mathbf{A} = \mathbf{U} \Lambda \mathbf{U}^{-1} = \mathbf{U} \Lambda \mathbf{U}^T \\ &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 0 & 0 & 1 \end{pmatrix} \end{aligned}$$

Matrices \mathbf{U} and Λ are much simpler than matrix \mathbf{A} . Moreover, \mathbf{U} and Λ are orthogonal and diagonal matrices and so they have many valuable properties.

$$tr(\Lambda) = \sum_{i=1}^n \lambda_i$$

$$rank(\Lambda) = n$$

$$\mathbf{U}^{-1} = \mathbf{U}^T$$

$$\mathbf{U}^{-1} \mathbf{U} = \mathbf{U} \mathbf{U}^{-1} = \mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_n$$

$$|\mathbf{A}| = |\Lambda| = \prod_{i=1}^n \lambda_i$$

$$|\mathbf{A}|^k = |\Lambda|^k = \prod_{i=1}^n \lambda_i^k$$

$$\Lambda^k = \text{diag}(\lambda_1^k, \lambda_2^k, \dots, \lambda_n^k) = \begin{pmatrix} \lambda_1^k & 0 & \dots & 0 \\ 0 & \lambda_2^k & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n^k \end{pmatrix}$$

$$A^k = U\Lambda^k U^{-1} = U\Lambda^k U^T = \sum_{i=1}^n \lambda_i^k u_i u_i^T$$

Singular value decomposition

When matrix $A(m,n)$ is not square matrix, there is a technique called singular value decomposition (SVD), a generalization of Jordan decomposition, which is used to decompose matrix $A(m,n)$. Let $\Lambda = \text{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_r^{1/2})$ be the eigenvalue matrix where λ_i is eigenvalue of AA^T and A^TA . Note that λ_i (s) are in descending ordering, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r$. In other words, r eigenvalues λ_i (s) are solutions of one of two following equations:

$$|AA^T - \lambda I_m| = 0$$

$$|A^TA - \lambda I_n| = 0$$

The number of eigenvalues is $r = \text{rank}(AA^T) = \text{rank}(A^TA) \leq \min(m, n)$. There are r eigenvectors u_i (s) corresponding to r eigenvalues, which are solutions of equation:

$$(AA^T - \lambda_i I_m)x = 0$$

There are r eigenvectors v_i (s) corresponding to r eigenvalues, which are solutions of equation:

$$(A^TA - \lambda_i I_n)x = 0$$

Let $U(m,r)$ and $V(n,r)$ be eigenvector matrices of r eigenvectors u_i (s) and r eigenvectors v_i (s). It is easy to infer that columns of U and V are mutually orthogonal eigenvectors and each eigenvector u_i (v_i) has m (n) components. The *singular value decomposition* (SVD) theorem states that (Härdle & Simar, 2013, p. 64):

$$A = U\Lambda V^T$$

Where,

$$U_{m \times r} = (u_1, u_2, \dots, u_r) = \begin{pmatrix} u_{11} & u_{21} & \dots & u_{r1} \\ u_{12} & u_{22} & \dots & u_{r2} \\ \vdots & \vdots & \ddots & \vdots \\ u_{1m} & u_{2m} & \dots & u_{rm} \end{pmatrix}$$

where $u_i = \begin{pmatrix} u_{i1} \\ u_{i2} \\ \vdots \\ u_{in} \end{pmatrix}$ is a column normalized eigenvector

$$V_{n \times r} = (v_1, v_2, \dots, v_r) = \begin{pmatrix} v_{11} & v_{21} & \dots & v_{r1} \\ v_{12} & v_{22} & \dots & v_{r2} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1n} & v_{2n} & \dots & v_{rn} \end{pmatrix}$$

where $\mathbf{v}_i = \begin{pmatrix} v_{i1} \\ v_{i2} \\ \vdots \\ v_{in} \end{pmatrix}$ is a column normalized eigenvector

$$\Lambda_{n \times r} = \begin{pmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda_r} \end{pmatrix}$$

According to SVD theorem, any matrix can be decomposed into three simpler matrices. SVD is often used to reduce vector space. For instance, when you choose $k \ll r$, which means that k is much smaller than the number r of eigenvalues, the number of columns of \mathbf{U} , Λ and \mathbf{V} are reduced to be $k \ll r$. Let $\mathbf{A}'(k \times k)$, $\mathbf{U}'(m \times k)$ and $\mathbf{V}'(n \times k)$ are eigenvalue matrix and eigenvector matrices with respect to k eigenvalues, we have:

$$\mathbf{A}' = \mathbf{U}' \mathbf{A}' \mathbf{V}'^T$$

Where $\mathbf{A}'(m \times n)$ is approximated matrix of \mathbf{A} .

In common, given matrix $\mathbf{A}(m \times n)$ in n -dimensional vector space with suppose that $m < n$ and it is necessary to reduce the dimension of vector space as small as possible. It means that the number k of descending-order eigenvalues ($\lambda_1 > \lambda_2 > \dots > \lambda_k$) is as small as possible and so spectrums of \mathbf{A} such as \mathbf{U}' , Λ' and \mathbf{V}' are very small while re-created matrix \mathbf{A}' is approximated to \mathbf{A} . In other words, the dimension of mentioned vector space is reduced to k .

In practical, there is another version of SVD (Baker, 2013, pp. 14-23). Without loss of generality, suppose $m < n$ and let $r \leq n$ be the number of eigenvalues resulting from two equations $|\mathbf{AA}^T - \lambda \mathbf{I}_m| = 0$ and $|\mathbf{AA}^T - \lambda \mathbf{I}_n|$. If $r < n$, it is possible to set $n - r$ eigenvalues to be zero and so we suppose that $r = n$ without loss of generality. There are m normalized eigenvectors \mathbf{u}_i (s) resulted from $(\mathbf{AA}^T - \lambda_i \mathbf{I}_m)\mathbf{x} = \mathbf{0}$ and n normalized eigenvectors \mathbf{v}_i (s) resulted from $(\mathbf{A}^T \mathbf{A} - \lambda_i \mathbf{I}_n)\mathbf{x} = \mathbf{0}$ and so two eigenmatrices \mathbf{U} and \mathbf{V} are totally determined as below.

$$\mathbf{U}_{m \times m} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m) = \begin{pmatrix} u_{11} & u_{21} & \dots & u_{m1} \\ u_{12} & u_{22} & \dots & u_{r2} \\ \vdots & \vdots & \ddots & \vdots \\ u_{1m} & u_{2m} & \dots & u_{mm} \end{pmatrix}$$

where $\mathbf{u}_i = \begin{pmatrix} u_{i1} \\ u_{i2} \\ \vdots \\ u_{im} \end{pmatrix}$ is a column normalized eigenvector

$$\mathbf{V}_{n \times n} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n) = \begin{pmatrix} v_{11} & v_{21} & \dots & v_{n1} \\ v_{12} & v_{22} & \dots & v_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1n} & v_{2n} & \dots & v_{nn} \end{pmatrix}$$

where $\mathbf{v}_i = \begin{pmatrix} v_{i1} \\ v_{i2} \\ \vdots \\ v_{in} \end{pmatrix}$ is a column normalized eigenvector

$$\Lambda_{m \times n} = \begin{pmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \dots & 0 & 0 & \dots & 0 \\ 0 & 0 & \ddots & 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \sqrt{\lambda_m} & 0 & \dots & 0 \end{pmatrix}$$

Please pay attention that $\mathbf{U}(m \times m)$ and $\mathbf{V}(n \times n)$ are orthogonal matrices because they are square matrices composed of orthogonal eigenvectors. Moreover eigenvalue matrix $\Lambda(m \times n)$ is not square matrix and it uses only m eigenvalues among n eigenvalues in its subsidiary diagonal and remaining elements are zeros. Matrix \mathbf{A} is decomposed by these matrices.

$$\mathbf{A} = \mathbf{U}\Lambda\mathbf{V}^T = \mathbf{U}\Lambda\mathbf{V}'^T$$

In this version, although eigenvalue matrix Λ is not diagonal, $\mathbf{U}(m \times m)$ and $\mathbf{V}(n \times n)$ are orthogonal matrices and so this gives us many valuable properties.

By reducing technique, it is possible to decrease the number of eigenvectors in matrix \mathbf{V} from n to m , which makes matrix \mathbf{V} become $m \times m$ orthogonal matrix and so Λ is $m \times m$ diagonal matrix. Because \mathbf{A} is invertible in this situation, let $\mathbf{A}^{-1} = \mathbf{V}\Lambda^{-1}\mathbf{U}^T$ and we have:

$$\mathbf{A}\mathbf{A}^{-1}\mathbf{A} = (\mathbf{U}\Lambda\mathbf{V}^T)(\mathbf{V}\Lambda^{-1}\mathbf{U}^T)(\mathbf{U}\Lambda\mathbf{V}^T) = \mathbf{U}\Lambda\mathbf{V}^T = \mathbf{A}$$

We infer that $\mathbf{A}\mathbf{A}^{-1}\mathbf{A}$ is G-inverse of \mathbf{A} .

Moreover, if you want to reduce vector space as small as possible, you can choose k descending-order eigenvalues ($\lambda_1 > \lambda_2 > \dots > \lambda_k$) with $k < m$ and $k \ll n$, which means that k is smaller than m and much smaller than n . So reduced eigenvalue and eigenvector matrices \mathbf{U}' , Λ' and \mathbf{V}' become very small.

$$\mathbf{U}'_{m \times k} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k) = \begin{pmatrix} u_{11} & u_{21} & \dots & u_{k1} \\ u_{12} & u_{22} & \dots & u_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ u_{1m} & u_{2m} & \dots & u_{km} \end{pmatrix}$$

where $\mathbf{u}_i = \begin{pmatrix} u_{i1} \\ u_{i2} \\ \vdots \\ u_{im} \end{pmatrix}$ is a column normalized eigenvector

$$\mathbf{U}'_{n \times k} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k) = \begin{pmatrix} v_{11} & v_{21} & \dots & v_{k1} \\ v_{12} & v_{22} & \dots & v_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1n} & v_{2n} & \dots & v_{kn} \end{pmatrix}$$

where $\mathbf{v}_i = \begin{pmatrix} v_{i1} \\ v_{i2} \\ \vdots \\ v_{in} \end{pmatrix}$ is a column normalized eigenvector

$$\Lambda'_{k,k} = \begin{pmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda_k} \end{pmatrix}$$

Note that re-created $\mathbf{A}' = \mathbf{U}'\mathbf{A}'\mathbf{V}'^T$ is approximated to matrix \mathbf{A} . Author (Baker, 2013, p. 18) gives an example for illustrating SVD with matrix $\mathbf{A} = \begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix}$. In order to find out eigenvalues and eigenvectors that compose eigenmatrix \mathbf{U} , we solve following equation.

$$|\mathbf{A}\mathbf{A}^T - \lambda I_2| = 0 \Leftrightarrow \left| \begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} \begin{pmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right| = 0 \Rightarrow \begin{cases} \lambda_1 = 12 \\ \lambda_2 = 10 \end{cases}$$

The eigenvectors \mathbf{u}_i creating eigenmatrix \mathbf{U} is found out via $\lambda_1 = 12$ and $\lambda_2 = 10$ as follows:

$$(\mathbf{A}\mathbf{A}^T - 12I_2)\mathbf{x} = \mathbf{0} \Leftrightarrow \left(\begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} \begin{pmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} - 12 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = x_1 \begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow \mathbf{u}_1 = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}$$

$$(\mathbf{A}\mathbf{A}^T - 10I_2)\mathbf{x} = \mathbf{0} \Leftrightarrow \left(\begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} \begin{pmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} - 10 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = x_1 \begin{pmatrix} 1 \\ -1 \end{pmatrix} \Rightarrow \mathbf{u}_2 = \begin{pmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{pmatrix}$$

So eigenmatrix \mathbf{U} is totally determined, $\mathbf{U} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix}$. Now we calculate second eigenmatrix \mathbf{V} by solving following equation.

$$|\mathbf{A}^T\mathbf{A} - \lambda I_2| = 0 \Leftrightarrow \left| \begin{pmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right| = 0 \Rightarrow \begin{cases} \lambda_1 = 12 \\ \lambda_2 = 10 \\ \lambda_3 = 0 \end{cases}$$

The eigenvectors \mathbf{v}_i composing eigenmatrix \mathbf{V} is found out via $\lambda_1 = 12$, $\lambda_2 = 10$ and $\lambda_3 = 0$ as follows:

$$(\mathbf{A}^T\mathbf{A} - 12I_2)\mathbf{x} = \mathbf{0} \Leftrightarrow \left(\begin{pmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} - 12 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = x_3 \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix} \Rightarrow \mathbf{v}_1 = \begin{pmatrix} 1/\sqrt{6} \\ 2/\sqrt{6} \\ 1/\sqrt{6} \end{pmatrix}$$

$$(\mathbf{A}^T \mathbf{A} - 10)\mathbf{x} = \mathbf{0} \Leftrightarrow \left(\begin{pmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} - 10 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = x_2 \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} \Rightarrow \mathbf{v}_2 = \begin{pmatrix} 2/\sqrt{5} \\ -1/\sqrt{5} \\ 0 \end{pmatrix}$$

$$(\mathbf{A}^T \mathbf{A} - 0)\mathbf{x} = \mathbf{0} \Leftrightarrow \left(\begin{pmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} - 0 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = \begin{pmatrix} 1 \\ 2 \\ -5 \end{pmatrix} \Rightarrow \mathbf{v}_3 = \begin{pmatrix} 1/\sqrt{30} \\ 2/\sqrt{30} \\ -5/\sqrt{30} \end{pmatrix}$$

So eigenmatrix \mathbf{V} is totally determined, in general, we have eigenvalue matrix Λ and eigenvector matrices \mathbf{U} and \mathbf{V} as follows:

$$\mathbf{U} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \text{ and } \Lambda = \begin{pmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{pmatrix}$$

$$\text{and } \mathbf{V} = \begin{pmatrix} 1/\sqrt{6} & 2/\sqrt{5} & 1/\sqrt{30} \\ 2/\sqrt{6} & -1/\sqrt{5} & 2/\sqrt{30} \\ 1/\sqrt{6} & 0 & -5/\sqrt{30} \end{pmatrix}$$

It is easy to validate spectrum decomposition of matrix $\mathbf{A}(m,n)$.

$$\begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} = \mathbf{A} = \mathbf{U} \Lambda \mathbf{V}^T$$

$$= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{pmatrix} \begin{pmatrix} 1/\sqrt{6} & 2/\sqrt{5} & 1/\sqrt{30} \\ 2/\sqrt{6} & -1/\sqrt{5} & 2/\sqrt{30} \\ 1/\sqrt{6} & 0 & -5/\sqrt{30} \end{pmatrix}$$

When reducing matrix \mathbf{A} , let $k = 2$ be the smaller number of eigenvalues while the total number of eigenvalues are 3, we have:

$$\mathbf{U}' = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \text{ and } \Lambda' = \begin{pmatrix} \sqrt{12} & 0 \\ 0 & \sqrt{10} \end{pmatrix} \text{ and } \mathbf{V}' = \begin{pmatrix} 1/\sqrt{6} & 2/\sqrt{5} \\ 2/\sqrt{6} & -1/\sqrt{5} \\ 1/\sqrt{6} & 0 \end{pmatrix}$$

And re-created matrix \mathbf{A}' is:

$$\mathbf{A}' = \mathbf{U}' \Lambda' \mathbf{V}^T = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} \sqrt{12} & 0 \\ 0 & \sqrt{10} \end{pmatrix} \begin{pmatrix} 1/\sqrt{6} & 2/\sqrt{6} & 1/\sqrt{6} \\ 2/\sqrt{5} & -1/\sqrt{5} & 0 \end{pmatrix} = \begin{pmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{pmatrix} = \mathbf{A}$$

If $k = 1$ and we have:

$$\mathbf{U}' = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix} \text{ and } \Lambda' = (\sqrt{12}) \text{ and } \mathbf{V}' = \begin{pmatrix} 1/\sqrt{6} \\ 2/\sqrt{6} \\ 1/\sqrt{6} \end{pmatrix}$$

The re-created matrix \mathbf{A}' becomes:

$$\mathbf{A}' = \mathbf{U}' \Lambda' \mathbf{V}'^T = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix} (\sqrt{12}) (1/\sqrt{6} \quad 2/\sqrt{6} \quad 1/\sqrt{6}) = \begin{pmatrix} 1 & 2 & 1 \\ 1 & 2 & 1 \end{pmatrix} \approx \mathbf{A}$$

Quadratic form

Another important subject in matrix analysis is quadratic form (Härdle & Simar, 2013, pp. 65-67), which is described in general way. Given a invertible matrix $\mathbf{A}(n,n)$ and a random vector \mathbf{x} , the quadratic form of \mathbf{x} denoted $Q(\mathbf{x})$ is:

$$Q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

Quadratic form is *positive definite* or *positive semi-definite* if $Q(\mathbf{x}) > 0$ or $Q(\mathbf{x}) \geq 0$, respectively. Matrix \mathbf{A} is called positive definite denoted $\mathbf{A} > 0$ or positive semi-definite denoted $\mathbf{A} \geq 0$ if quadratic form $Q(\mathbf{x})$ of any vector with regard to \mathbf{A} is positive definite or positive semi-definite, respectively. Otherwise, if $Q(\mathbf{x}) < 0$ then, quadratic form and matrix \mathbf{A} are indefinite. Because \mathbf{A} is invertible, it is decomposed into $\mathbf{A} = \mathbf{U} \Lambda \mathbf{U}^T$ where Λ and \mathbf{U} are eigenvalue matrix and eigenvector matrix. Let vector $\mathbf{y} = \mathbf{U}^T \mathbf{x}$, we have (Härdle & Simar, 2013, p. 65):

$$Q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T (\mathbf{U} \Lambda \mathbf{U}^T) \mathbf{x} = \mathbf{y}^T \Lambda \mathbf{y} = \sum_{i=1}^n \lambda_i y_i^2$$

Where λ_i (s) are eigenvalues.

We have some properties of quadratic form:

$\mathbf{A} > 0$ if and only if all eigenvalues λ_i (s) > 0

If $\mathbf{A} > 0$ then \mathbf{A}^{-1} exists and $|\mathbf{A}| > 0$

Given invertible and symmetric matrices \mathbf{A} and \mathbf{B} and given a constraint $\mathbf{x}^T \mathbf{B} \mathbf{x} = 1$, the maximum (minimum) of $\mathbf{x}^T \mathbf{A} \mathbf{x}$ is the largest (smallest) eigenvalue of $\mathbf{B}^{-1} \mathbf{A}$. The vector maximizing (minimizing) $\mathbf{x}^T \mathbf{A} \mathbf{x}$ is the eigenvector which corresponds to largest (smallest) eigenvalue of $\mathbf{B}^{-1} \mathbf{A}$. We have (Härdle & Simar, 2013, p. 66):

$$\max_{x^T B x = 1} (x^T A x) = \max(\lambda_1, \lambda_2, \dots, \lambda_n)$$

$$\min_{x^T B x = 1} (x^T A x) = \min(\lambda_1, \lambda_2, \dots, \lambda_n)$$

Partitioned matrix

Now we research partitioned matrix (Härdle & Simar, 2013, pp. 68-70) which is useful in matrix analysis. Given matrix $A(m,n)$ are composed of other matrices $A_{ij}(m_i \times n_j)$.

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$$

Where A_{ij} (s) are matrices having m_i rows and n_j columns. A_{ij} (s) are called groups or sub-matrices and A is called partitioned matrix. Partitioning matrix technique performs matrix operations and determine properties and characteristics of matrix (determinant $|A|$, inverse A^{-1} , etc.) according to sub-matrices. Given partitioned matrix $B(m,n) = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$ similar to A , we have (Härdle & Simar, 2013, p. 69):

$$\begin{aligned} A + B &= \begin{pmatrix} A_{11} + B_{11} & A_{12} + B_{12} \\ A_{21} + B_{21} & A_{22} + B_{22} \end{pmatrix} \\ B^T &= \begin{pmatrix} B_{11}^T & B_{21}^T \\ B_{12}^T & B_{22}^T \end{pmatrix} \\ AB^T &= \begin{pmatrix} A_{11}B_{11}^T + A_{12}B_{12}^T & A_{11}B_{21}^T + A_{12}B_{22}^T \\ A_{21}B_{11}^T + A_{22}B_{12}^T & A_{21}B_{21}^T + A_{22}B_{22}^T \end{pmatrix} \end{aligned}$$

If A is invertible, the inverse of A is (Härdle & Simar, 2013, p. 69):

$$A^{-1} = \begin{pmatrix} A^{11} & A^{12} \\ A^{21} & A^{22} \end{pmatrix}$$

$$\text{Where } \begin{cases} A^{11} = (A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1} \\ A^{12} = -A^{11}A_{12}A_{22}^{-1} \\ A^{21} = -A_{22}^{-1}A_{21}A^{11} \\ A^{22} = A_{22}^{-1} + A_{22}^{-1}A_{21}A^{11}A_{12}A_{22}^{-1} \end{cases}$$

If A_{11} is invertible, the determinant of A is:

$$|A| = |A_{11}| |A_{22} - A_{21}A_{11}^{-1}A_{12}|$$

If A_{22} is invertible, the determinant of A is:

$$|A| = |A_{22}| |A_{11} - A_{12}A_{22}^{-1}A_{21}|$$

Suppose partitioned matrix B is composed of a invertible matrix $A(n,n)$ and two ($n, 1$) vectors a and b as follows:

$$B = \begin{pmatrix} 1 & b^T \\ a & A \end{pmatrix}$$

We have (Härdle & Simar, 2013, p. 69):

$$|\mathbf{B}| = |\mathbf{A} - \mathbf{ab}^T| = |\mathbf{A}| |1 - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{a}|$$

$$(\mathbf{A} - \mathbf{ab}^T)^{-1} = \mathbf{A}^{-1} + \frac{\mathbf{A}^{-1} \mathbf{ab}^T \mathbf{A}^{-1}}{1 - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{a}}$$

3. Matrix derivative

Recall that matrix analysis and matrix calculus are typical applications of matrix algebra. Matrix analysis focuses on processing and analyzing multivariate data while matrix calculus focuses on derivative and differential with regard to matrix. Matrix derivative is the basic concept of matrix calculus and differential, which extends principles of derivative from real number space \mathbb{R}^n to vector space. Vector space is a general space in which vectors and matrices are elements in vector space. Real number or scalar is 1-component vector (matrix). This section describes and classifies matrix derivatives. Its main contents are extracted from the webpage (Wikipedia, Matrix calculus, 2014) which is a valued document that readers should read. The term “scalar” refers to a real number but the study can be extended to complex number. Given a function f from n -dimension domain space \mathbb{R}^n to real image space R , in other words, domain of f is vector space and image of f is scalar space. For convenience, we have convention that vector space \mathbb{R}^n is based on real number field \mathbb{R} and default vector space is \mathbb{R}^n having n dimensions if there is no note.

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Thus, f is a scalar function with respect to (w.r.t) vector variable \mathbf{x} . Note that the domain space \mathbb{R}^n can be reduced to scalar space \mathbb{R} and f becomes scalar function over scalar space. The row derivative of f with respect to \mathbf{x} is defined as below:

$$\frac{\partial f}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

So the row derivative of f w.r.t \mathbf{x} is a row vector whose components are partial derivatives $\frac{\partial f}{\partial x_i}$ where x_i (s) are components of \mathbf{x} . Similarly, the column derivative of f w.r.t \mathbf{x} is:

$$\frac{\partial f}{\partial \mathbf{x}^T} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

The derivative of scalar function f w.r.t vector variable \mathbf{x} is often called *gradient* of f denoted $gradf$ or ∇f . Gradient ∇f is often known as row derivative vector. In general, we have seven ways to denote derivative as following.

$$f'(\mathbf{x}) = gradf = \nabla f = Df = \frac{df}{d\mathbf{x}} = \frac{\partial f}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial \mathbf{x}^T} \right)^T$$

Df and $\frac{df}{d\mathbf{x}}$ are formal notations originated from differential operation which is only applied into real number field \mathbb{R} and vector space \mathbb{R}^n . When derivative gets value at \mathbf{x}_0 , we use following notations with the same meaning:

$$f'(x_0) = \text{grad}f(x_0) = \nabla f(x_0) = Df(x_0) = \frac{df}{dx}(x_0) = \frac{\partial f}{\partial x}(x_0) = \left(\frac{\partial f}{\partial x^T}\right)^T (x_0)$$

The notation $f'(\mathbf{x})$ is the most popular in number analysis but we prefer to use notation $\frac{\partial f}{\partial \mathbf{x}}$ and $\frac{\partial f}{\partial \mathbf{x}^T}$ in matrix calculus. When the image space is extended to vector space \mathbb{R}^m , we have vector function \mathbf{f} .

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} \mapsto \mathbf{f} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{pmatrix}$$

Each elemental function $f_i \in \mathbf{f}$ may be the scalar function of all x_j (s). The derivative of vector function \mathbf{f} with respective to vector variable \mathbf{x} is a so-called Jacobian matrix.

$$\mathbf{f}'(\mathbf{x}) = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

Where $\frac{\partial f_i}{\partial x_j}$ is the partial derivative of i^{th} functional component f_i with respect to j^{th} domain component x_j . Aforementioned scalar function f and vector function \mathbf{f} are typical examples for matrix derivative. If variable \mathbf{x} in domain space is considered independent variable and function f in image space is considered dependent variable. We follow the same convention that “normal letters denote scalar; vectors are denoted by bold letters and matrices are denoted as bold and uppercase letters”. There are nine possible combinations (Wikipedia, Matrix calculus, 2014) between independent variable (function) and dependent variable, each of combination corresponds a kind of matrix derivative when scalar and vector are simple forms of matrix. Concretely, scalar is 1×1 matrix, column vector is $n \times 1$ matrix and row vector is $1 \times n$ matrix.

Dependent variable (function)	Independent variable		
	Scalar	Vector	Matrix
Scalar	$\frac{\partial f}{\partial x}$	$\frac{\partial f}{\partial \mathbf{x}}$	$\frac{\partial f}{\partial \mathbf{X}}$
Vector	$\frac{\partial \mathbf{f}}{\partial x}$	$\frac{\partial \mathbf{f}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{f}}{\partial \mathbf{X}}$
Matrix	$\frac{\partial \mathbf{F}}{\partial x}$	$\frac{\partial \mathbf{F}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{F}}{\partial \mathbf{X}}$

In table (Wikipedia, Matrix calculus, 2014) above, there are six descriptive kinds of derivative: scalar-by-scalar (scalar function w.r.t scalar variable), scalar-by-vector (scalar function w.r.t vector variable), scalar-by-matrix (scalar function w.r.t matrix variable), vector-by-scalar (vector function

w.r.t scalar variable), vector-by-vector (vector function w.r.t vector variable) and matrix-by-scalar (matrix function w.r.t vector variable). Three remaining kinds corresponding to shaded cells such as vector-by-matrix, matrix-by-vector and matrix-by-matrix, which relate to tensor product, are not focused in the report. We divide these kinds of function into three groups:

- Scalar function group includes scalar-by-scalar, scalar-by-vector and scalar-by-matrix.
- Vector function group includes vector-by-scalar, vector-by-vector and vector-by-matrix.
- Matrix function group includes matrix-by-scalar, matrix-by-vector and matrix-by-matrix.

Suppose that all mentioned functions (dependent variables) belong to differentiability class, which means that their k^{th} derivatives exist and continuous. Before discussing matrix derivative in more detailed, it is convenient for us to apply *numerator layout convention* (Wikipedia, Matrix calculus, 2014) into denoting multivariate derivatives. With numerator layout convention, given the vector-by-vector function from \mathbb{R}^n to \mathbb{R}^m , the derivative $\frac{\partial \mathbf{f}}{\partial \mathbf{x}}$ is a Jacobian matrix in which the numerator \mathbf{f}

(dependent variable) lists its partial derivatives $\begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{pmatrix}$ as column vector and the independent variable \mathbf{x} lists its components as row vector (x_1, x_2, \dots, x_n) . In other words, numerator layout convention results out partial derivatives according to column alignment with \mathbf{f} and row alignment with \mathbf{x}^T .

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

Let $\mathbf{y} = \mathbf{f}(\mathbf{x})$ as dependent variable, we have convention that \mathbf{y} is identical to \mathbf{f} . Please pay attention to this convention because we will often use dependent variable \mathbf{y} instead of \mathbf{f} .

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

Now we discuss more about six kinds of derivative: scalar-by-scalar, scalar-by-vector, scalar-by-matrix, vector-by-scalar, vector-by-vector and matrix-by-scalar.

Scalar-by-scalar derivative

Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be scalar-by-scalar function, the derivative of scalar y by scalar x is $f'(x) = \frac{\partial y}{\partial x}$, which is the simplest case and so it is not mentioned much in matrix calculus. If $x = x(t)$ is scalar function of $t \in R$, the derivative of f with respect to variable t is $f'(t) = f'(x)x'(t) = \frac{\partial y}{\partial x} \frac{\partial x}{\partial t}$. This is chain rule which is applied into all kinds of derivative. The second-order derivative of f is $f''(x) = \frac{\partial^2 y}{\partial x^2}$. It is necessary to distinguish between second-order derivative and second power of derivative, $\frac{\partial^2 y}{\partial x^2} \neq \left(\frac{\partial y}{\partial x}\right)^2$.

Scalar-by-vector derivative

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be scalar-by-vector function, the derivative of scalar y by vector $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ is:

$$\frac{\partial y}{\partial \mathbf{x}} = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

This derivative is also called *gradient* of f , denoted $\text{grad } f$ or ∇f . Scalar-by-vector function is often called multivariate function. The *directional derivative* of scalar-by-vector function is dot product

between gradient ∇f and directional vector $\mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}$, which is denoted $\nabla_{\mathbf{u}} f$.

$$\nabla_{\mathbf{u}} f = \nabla f(\mathbf{x}) \mathbf{u} = \frac{\partial y}{\partial \mathbf{x}} \mathbf{u} = \left(\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_n} \right) \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}$$

The second derivative of scalar-by-vector function is a so-called Hessian matrix whose elements are partial second-order derivatives w.r.t partial variables.

$$f''(\mathbf{x}) = \nabla^2 f = \frac{\partial^2 y}{\partial \mathbf{x}^2} = \begin{pmatrix} \frac{\partial^2 y}{\partial x_1^2} & \frac{\partial^2 y}{\partial x_1 x_2} & \cdots & \frac{\partial^2 y}{\partial x_1 x_n} \\ \frac{\partial^2 y}{\partial x_2 x_1} & \frac{\partial^2 y}{\partial x_2^2} & \cdots & \frac{\partial^2 y}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 y}{\partial x_n x_1} & \frac{\partial^2 y}{\partial x_n x_2} & \cdots & \frac{\partial^2 y}{\partial x_n^2} \end{pmatrix}$$

Where $\frac{\partial^2 y}{\partial x_i^2}$ is the second derivative w.r.t partial variable x_i and $\frac{\partial^2 y}{\partial x_i x_j} = \frac{\partial}{\partial x_j} \left(\frac{\partial y}{\partial x_i} \right)$ is the iterative derivative that y is taken derivative with respect to x_i , which in turn, is taken with respect to x_j .

Scalar-by-matrix derivative

Let $f: \mathbb{R}^{n_x m} \rightarrow \mathbb{R}$ be scalar-by-matrix function, the derivative of scalar y by matrix $\mathbf{X}(m,n) =$

$$\begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix}$$

$$\frac{\partial y}{\partial \mathbf{X}} = \nabla f = \begin{pmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{21}} & \dots & \frac{\partial y}{\partial x_{m1}} \\ \frac{\partial y}{\partial x_{12}} & \frac{\partial y}{\partial x_{22}} & \dots & \frac{\partial y}{\partial x_{m2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{1n}} & \frac{\partial y}{\partial x_{2n}} & \dots & \frac{\partial y}{\partial x_{mn}} \end{pmatrix}$$

Where $\frac{\partial y}{\partial x_{ij}}$ is the partial derivative of y w.r.t partial element x_{ij} of matrix variable \mathbf{X} . While \mathbf{X} is $m \times n$ matrix, the derivative $\frac{\partial y}{\partial \mathbf{X}}$ is $n \times m$ matrix. This follows the numerator layout convention and so element indices in derivative matrix $\frac{\partial y}{\partial \mathbf{X}}$ are transposed as compared with element indices in variable \mathbf{X} . In similar to scalar-by-vector function, the scalar-by-matrix derivative is also called *gradient matrix* ∇f of f . The directional derivative of scalar y by matrix $\mathbf{X}(m,n)$ in the direction of matrix $\mathbf{U}(m,n)$ is defined via the trace of matrix mentioned in previous section. So the directional derivative denoted $\nabla_{\mathbf{U}} f$ is a scalar and its formula is:

$$\nabla_{\mathbf{U}} f = \nabla f(\mathbf{X}) \mathbf{U} = \text{tr} \left(\frac{\partial y}{\partial \mathbf{X}} \mathbf{U} \right)$$

Note that the above formula is also true to directional derivative of scalar-by-vector function as aforementioned when the trace of scalar is itself. We conclude that concept “gradient” always goes along with concept “directional derivative” in scalar function and the evaluation of directional derivative is always scalar.

Vector-by-scalar derivative

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ be vector-by-scalar function, the derivative of vector $\mathbf{y} = f(x) = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$ by scalar x is:

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{pmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_n}{\partial x} \end{pmatrix}$$

Where $\frac{\partial y_i}{\partial x}$ is the derivative of partial element y_i w.r.t scalar x . When f is vector function, its derivative $\frac{\partial \mathbf{y}}{\partial x}$ is known as *tangent vector*. Vector function is applied into representing curve and surface in differential geometry and tangent vector is important component w.r.t geometrical object. The second-order derivative of vector-by-scalar function is:

$$\frac{\partial^2 \mathbf{y}}{\partial x^2} = \begin{pmatrix} \frac{\partial^2 y_1}{\partial x^2} \\ \frac{\partial^2 y_2}{\partial x^2} \\ \vdots \\ \frac{\partial^2 y_n}{\partial x^2} \end{pmatrix}$$

Vector-by-vector derivative

Given another vector function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is vector-by-vector function from space \mathbb{R}^n to space \mathbb{R}^m .

The derivative of vector $\mathbf{y} = f(\mathbf{x}) = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$ by scalar $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix}$ is:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

Where $\frac{\partial y_i}{\partial x_j}$ is the partial derivative of i^{th} functional component y_i with respect to j^{th} domain component x_j . The derivative matrix as aforementioned is known Jacobian matrix, push-forward matrix, or differential matrix. Jacobian matrix is the heart of matrix calculus. If we consider vector function \mathbf{y} is the composition of partial scalar function y_1, y_2, \dots, y_m then Jacobian matrix is the column vector of gradients $\nabla f_i = \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}}$.

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \nabla f_1 \\ \nabla f_2 \\ \vdots \\ \nabla f_m \end{pmatrix} = \begin{pmatrix} \frac{\partial \mathbf{y}_1}{\partial \mathbf{x}} \\ \frac{\partial \mathbf{y}_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial \mathbf{y}_m}{\partial \mathbf{x}} \end{pmatrix} \text{ where gradient } \nabla f_i = \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}} = \left(\frac{\partial y_i}{\partial x_1}, \frac{\partial y_i}{\partial x_2}, \dots, \frac{\partial y_i}{\partial x_n} \right)$$

Matrix-by-scalar derivative

Let $\mathbf{F}: \mathbb{R} \rightarrow \mathbb{R}^{n \times m}$ be matrix-by-scalar function, the derivative of matrix $\mathbf{Y} = \mathbf{F}(x) =$

$$\begin{pmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ y_{21} & y_{22} & \dots & y_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ y_{m1} & y_{m2} & \dots & y_{mn} \end{pmatrix}$$

by scalar x is:

$$\frac{\partial \mathbf{Y}}{\partial x} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x} & \frac{\partial y_{12}}{\partial x} & \dots & \frac{\partial y_{1n}}{\partial x} \\ \frac{\partial y_{21}}{\partial x} & \frac{\partial y_{22}}{\partial x} & \dots & \frac{\partial y_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{m1}}{\partial x} & \frac{\partial y_{m2}}{\partial x} & \dots & \frac{\partial y_{mn}}{\partial x} \end{pmatrix}$$

Where $\frac{\partial y_{ij}}{\partial x}$ is the derivative of partial element y_{ij} w.r.t scalar x . When \mathbf{F} is matrix function, its derivative $\frac{\partial \mathbf{F}}{\partial x}$ is known as *tangent matrix*. The second-order derivative of matrix-by-scalar function is:

$$\frac{\partial^2 \mathbf{Y}}{\partial x^2} = \begin{pmatrix} \frac{\partial^2 y_{11}}{\partial x^2} & \frac{\partial^2 y_{12}}{\partial x^2} & \dots & \frac{\partial^2 y_{1n}}{\partial x^2} \\ \frac{\partial^2 y_{21}}{\partial x^2} & \frac{\partial^2 y_{22}}{\partial x^2} & \dots & \frac{\partial^2 y_{2n}}{\partial x^2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 y_{m1}}{\partial x^2} & \frac{\partial^2 y_{m2}}{\partial x^2} & \dots & \frac{\partial^2 y_{mn}}{\partial x^2} \end{pmatrix}$$

The derivative of matrix inverse denoted $\frac{\partial \mathbf{Y}^{-1}}{\partial x}$ is calculated based on tangent matrix $\frac{\partial \mathbf{Y}}{\partial x}$:

$$\frac{\partial \mathbf{Y}^{-1}}{\partial x} = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \mathbf{Y}^{-1}$$

Now we researched six kinds of derivative such as scalar-by-scalar, scalar-by-vector, scalar-by-matrix, vector-by-scalar, vector-by-vector and matrix-by-scalar. There are three remaining kinds of derivative such as vector-by-matrix, matrix-by-vector and matrix-by-matrix not concerned much in the report. Moreover high-order derivatives such as second-order, third order and k^{th} order derivatives are mentioned in restriction because they relate to tensor product – a complicated subject which goes beyond this content. We will discuss tensor product, high-order derivative in detail in another report focusing on the subject of combination of tensor product and multidimensional derivative. However we should have an overview of some concepts of derivative relating to matrix such as *vector-by-matrix*, *matrix-by-vector* and *matrix-by-matrix*.

Vector-by-matrix derivative

Let $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$ be vector-by-matrix function from $\mathbb{R}^{m \times n}$ space to \mathbb{R}^p space, the derivative of vector $\mathbf{y} = f(\mathbf{X}) = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{pmatrix}$ by matrix $\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix}$ is defined as follows:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial \mathbf{y}}{\partial x_{11}} & \frac{\partial \mathbf{y}}{\partial x_{21}} & \dots & \frac{\partial \mathbf{y}}{\partial x_{m1}} \\ \frac{\partial \mathbf{y}}{\partial x_{12}} & \frac{\partial \mathbf{y}}{\partial x_{22}} & \dots & \frac{\partial \mathbf{y}}{\partial x_{m2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{y}}{\partial x_{1n}} & \frac{\partial \mathbf{y}}{\partial x_{2n}} & \dots & \frac{\partial \mathbf{y}}{\partial x_{mn}} \end{pmatrix}$$

$$\text{where each element } \frac{\partial \mathbf{y}}{\partial x_{ij}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_{ij}} \\ \frac{\partial y_2}{\partial x_{ij}} \\ \vdots \\ \frac{\partial y_p}{\partial x_{ij}} \end{pmatrix} \text{ is a partial tangent vector}$$

Note that indices of matrix variable \mathbf{X} are transposed due to numerator layout convention. The interesting discovery is that the vector-by-matrix derivative $\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$ is compound gradient matrix whose components are partial tangent vectors.

Matrix-by-vector derivative

Let $F: \mathbb{R}^n \rightarrow \mathbb{R}^{p \times q}$ be matrix-by-vector function from \mathbb{R}^n space to $\mathbb{R}^{p \times q}$ space, the derivative of matrix $\mathbf{Y} = F(\mathbf{x}) = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1q} \\ y_{21} & y_{22} & \dots & y_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ y_{p1} & y_{p2} & \dots & y_{pq} \end{pmatrix}$ by vector $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ is defined as follows:

$$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}} = \left(\frac{\partial \mathbf{Y}}{\partial x_1}, \frac{\partial \mathbf{Y}}{\partial x_2}, \dots, \frac{\partial \mathbf{Y}}{\partial x_n} \right)$$

$$\text{where each element } \frac{\partial \mathbf{Y}}{\partial x_i} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x_i} & \frac{\partial y_{12}}{\partial x_i} & \dots & \frac{\partial y_{1q}}{\partial x_i} \\ \frac{\partial y_{21}}{\partial x_i} & \frac{\partial y_{22}}{\partial x_i} & \dots & \frac{\partial y_{2q}}{\partial x_i} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{p1}}{\partial x_i} & \frac{\partial y_{p2}}{\partial x_i} & \dots & \frac{\partial y_{pq}}{\partial x_i} \end{pmatrix} \text{ is a partial tangent matrix}$$

The interesting discovery is that the matrix-by-vector derivative $\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$ is compound gradient vector whose components are partial tangent matrices.

Matrix-by-matrix derivative

Given the most general case – matrix-by-matrix function $\mathbf{F}: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{p \times q}$ from $\mathbb{R}^{n \times m}$ space to $\mathbb{R}^{p \times q}$ space, the derivative of matrix

$$\mathbf{Y} = \mathbf{F}(\mathbf{X}) = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1q} \\ y_{21} & y_{22} & \dots & y_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ y_{p1} & y_{p2} & \dots & y_{pq} \end{pmatrix} \text{ by matrix } \mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix}$$

is defined as follows:

$$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial \mathbf{Y}}{\partial x_{11}} & \frac{\partial \mathbf{Y}}{\partial x_{21}} & \dots & \frac{\partial \mathbf{Y}}{\partial x_{m1}} \\ \frac{\partial \mathbf{Y}}{\partial x_{12}} & \frac{\partial \mathbf{Y}}{\partial x_{22}} & \dots & \frac{\partial \mathbf{Y}}{\partial x_{m2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{Y}}{\partial x_{1n}} & \frac{\partial \mathbf{Y}}{\partial x_{2n}} & \dots & \frac{\partial \mathbf{Y}}{\partial x_{mn}} \end{pmatrix} \text{ where } \frac{\partial \mathbf{Y}}{\partial x_{ij}} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x_{ij}} & \frac{\partial y_{12}}{\partial x_{ij}} & \dots & \frac{\partial y_{1q}}{\partial x_{ij}} \\ \frac{\partial y_{21}}{\partial x_{ij}} & \frac{\partial y_{22}}{\partial x_{ij}} & \dots & \frac{\partial y_{2q}}{\partial x_{ij}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{p1}}{\partial x_{ij}} & \frac{\partial y_{p2}}{\partial x_{ij}} & \dots & \frac{\partial y_{pq}}{\partial x_{ij}} \end{pmatrix}$$

The matrix-by-matrix derivative is $m \times n$ matrix whose elements are $p \times q$ matrices; in other words, the derivative is a fourth-rank tensor. The interesting discovery is that the matrix-by-matrix derivative $\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$ is compound gradient matrix whose components are partial tangent matrices.

In general, following is the summary table concerning derivatives divided into 3 groups such as scalar function, vector function and matrix function with 9 kinds of functions: scalar-by-scalar, scalar-by-vector, scalar-by-matrix, vector-by-scalar, vector-by-vector, vector-by-matrix, matrix-by-scalar, matrix-by-vector and matrix-by-matrix.

Dependent variable (function)	Independent variable		
	Scalar	Vector	Matrix
Scalar	Derivative is: $\frac{\partial y}{\partial x} = f'(x)$ Second-order derivative is: $\frac{\partial^2 y}{\partial x^2} = f''(x)$	Derivative known as gradient vector is: $\frac{\partial y}{\partial \mathbf{x}} = \nabla f$ $= \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$ Directional derivative is: $\frac{\partial y}{\partial \mathbf{x}} \mathbf{u} = \nabla_u f = \nabla f(\mathbf{x}) \mathbf{u} =$ $\left(\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_n} \right) \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}$ Second-order derivative known as Hessian matrix is:	Derivative known as gradient matrix is: $\frac{\partial \mathbf{y}}{\partial \mathbf{X}} = \nabla f =$ $\begin{pmatrix} \frac{\partial y_{11}}{\partial x_{11}} & \frac{\partial y_{12}}{\partial x_{12}} & \dots & \frac{\partial y_{1n}}{\partial x_{1n}} \\ \frac{\partial y_{21}}{\partial x_{21}} & \frac{\partial y_{22}}{\partial x_{22}} & \dots & \frac{\partial y_{2n}}{\partial x_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{p1}}{\partial x_{p1}} & \frac{\partial y_{p2}}{\partial x_{p2}} & \dots & \frac{\partial y_{pn}}{\partial x_{pn}} \end{pmatrix}$ Directional derivative is: $\nabla_u f = \nabla f(\mathbf{X}) \mathbf{U}$ $= \text{tr} \left(\frac{\partial \mathbf{y}}{\partial \mathbf{X}} \mathbf{U} \right)$

		$\frac{\partial^2 y}{\partial \mathbf{x}^2} = \nabla^2 f =$ $\begin{pmatrix} \frac{\partial^2 y}{\partial x_1^2} & \frac{\partial^2 y}{\partial x_1 x_2} & \cdots & \frac{\partial^2 y}{\partial x_1 x_n} \\ \frac{\partial^2 y}{\partial x_2 x_1} & \frac{\partial^2 y}{\partial x_2^2} & \cdots & \frac{\partial^2 y}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 y}{\partial x_n x_1} & \frac{\partial^2 y}{\partial x_n x_2} & \cdots & \frac{\partial^2 y}{\partial x_n^2} \end{pmatrix}$	
Vector	Derivative known as <i>tangent vector</i> is: $\frac{\partial \mathbf{y}}{\partial x} = \begin{pmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_n}{\partial x} \end{pmatrix}$ Second-order derivative is: $\frac{\partial^2 \mathbf{y}}{\partial x^2} = \begin{pmatrix} \frac{\partial^2 y_1}{\partial x^2} \\ \frac{\partial^2 y_2}{\partial x^2} \\ \vdots \\ \frac{\partial^2 y_n}{\partial x^2} \end{pmatrix}$	Derivative known as <i>Jacobian matrix</i> is: $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{pmatrix} =$ $\begin{pmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{pmatrix} = \begin{pmatrix} \nabla f_1 \\ \nabla f_2 \\ \vdots \\ \nabla f_m \end{pmatrix}$	Derivative known as compound gradient matrix is: $\frac{\partial \mathbf{y}}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial \mathbf{y}}{\partial x_{11}} & \frac{\partial \mathbf{y}}{\partial x_{21}} & \cdots & \frac{\partial \mathbf{y}}{\partial x_{m1}} \\ \frac{\partial \mathbf{y}}{\partial x_{12}} & \frac{\partial \mathbf{y}}{\partial x_{22}} & \cdots & \frac{\partial \mathbf{y}}{\partial x_{m2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{y}}{\partial x_{1n}} & \frac{\partial \mathbf{y}}{\partial x_{2n}} & \cdots & \frac{\partial \mathbf{y}}{\partial x_{mn}} \end{pmatrix}$ Where each element known as partial tangent vector is $\frac{\partial \mathbf{y}}{\partial x_{ij}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_{ij}} \\ \frac{\partial y_2}{\partial x_{ij}} \\ \vdots \\ \frac{\partial y_p}{\partial x_{ij}} \end{pmatrix}$
Matrix	Derivative known as <i>tangent matrix</i> is: $\frac{\partial \mathbf{Y}}{\partial x} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x} & \frac{\partial y_{12}}{\partial x} & \cdots & \frac{\partial y_{1n}}{\partial x} \\ \frac{\partial y_{21}}{\partial x} & \frac{\partial y_{22}}{\partial x} & \cdots & \frac{\partial y_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{m1}}{\partial x} & \frac{\partial y_{m2}}{\partial x} & \cdots & \frac{\partial y_{mn}}{\partial x} \end{pmatrix}$ Second-order derivative is: $\frac{\partial^2 \mathbf{Y}}{\partial x^2} = \begin{pmatrix} \frac{\partial^2 y_{11}}{\partial x^2} & \frac{\partial^2 y_{12}}{\partial x^2} & \cdots & \frac{\partial^2 y_{1n}}{\partial x^2} \\ \frac{\partial^2 y_{21}}{\partial x^2} & \frac{\partial^2 y_{22}}{\partial x^2} & \cdots & \frac{\partial^2 y_{2n}}{\partial x^2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 y_{m1}}{\partial x^2} & \frac{\partial^2 y_{m2}}{\partial x^2} & \cdots & \frac{\partial^2 y_{mn}}{\partial x^2} \end{pmatrix}$ Derivative of inverse is:	Derivative known as compound gradient vector is: $\frac{\partial \mathbf{Y}}{\partial \mathbf{x}} = \left(\frac{\partial \mathbf{Y}}{\partial x_1}, \frac{\partial \mathbf{Y}}{\partial x_2}, \dots, \frac{\partial \mathbf{Y}}{\partial x_n} \right)$ Where each element known as partial tangent matrix is: $\frac{\partial \mathbf{Y}}{\partial x_i} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x_i} & \frac{\partial y_{12}}{\partial x_i} & \cdots & \frac{\partial y_{1n}}{\partial x_i} \\ \frac{\partial y_{21}}{\partial x_i} & \frac{\partial y_{22}}{\partial x_i} & \cdots & \frac{\partial y_{2n}}{\partial x_i} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{p1}}{\partial x_i} & \frac{\partial y_{p2}}{\partial x_i} & \cdots & \frac{\partial y_{pq}}{\partial x_i} \end{pmatrix}$ $\frac{\partial \mathbf{Y}}{\partial x_{ij}} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x_{ij}} & \frac{\partial y_{12}}{\partial x_{ij}} & \cdots & \frac{\partial y_{1q}}{\partial x_{ij}} \\ \frac{\partial y_{21}}{\partial x_{ij}} & \frac{\partial y_{22}}{\partial x_{ij}} & \cdots & \frac{\partial y_{2q}}{\partial x_{ij}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{p1}}{\partial x_{ij}} & \frac{\partial y_{p2}}{\partial x_{ij}} & \cdots & \frac{\partial y_{pq}}{\partial x_{ij}} \end{pmatrix}$ Where each element known as partial tangent matrix is $\frac{\partial \mathbf{Y}}{\partial x_{ij}} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x_{ij}} & \frac{\partial y_{12}}{\partial x_{ij}} & \cdots & \frac{\partial y_{1q}}{\partial x_{ij}} \\ \frac{\partial y_{21}}{\partial x_{ij}} & \frac{\partial y_{22}}{\partial x_{ij}} & \cdots & \frac{\partial y_{2q}}{\partial x_{ij}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{p1}}{\partial x_{ij}} & \frac{\partial y_{p2}}{\partial x_{ij}} & \cdots & \frac{\partial y_{pq}}{\partial x_{ij}} \end{pmatrix}$	Derivative known as compound gradient matrix is: $\frac{\partial \mathbf{Y}}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial \mathbf{Y}}{\partial x_{11}} & \frac{\partial \mathbf{Y}}{\partial x_{21}} & \cdots & \frac{\partial \mathbf{Y}}{\partial x_{m1}} \\ \frac{\partial \mathbf{Y}}{\partial x_{12}} & \frac{\partial \mathbf{Y}}{\partial x_{22}} & \cdots & \frac{\partial \mathbf{Y}}{\partial x_{m2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{Y}}{\partial x_{1n}} & \frac{\partial \mathbf{Y}}{\partial x_{2n}} & \cdots & \frac{\partial \mathbf{Y}}{\partial x_{mn}} \end{pmatrix}$ Where each element known as partial tangent matrix is $\frac{\partial \mathbf{Y}}{\partial x_{ij}} = \begin{pmatrix} \frac{\partial y_{11}}{\partial x_{ij}} & \frac{\partial y_{12}}{\partial x_{ij}} & \cdots & \frac{\partial y_{1q}}{\partial x_{ij}} \\ \frac{\partial y_{21}}{\partial x_{ij}} & \frac{\partial y_{22}}{\partial x_{ij}} & \cdots & \frac{\partial y_{2q}}{\partial x_{ij}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{p1}}{\partial x_{ij}} & \frac{\partial y_{p2}}{\partial x_{ij}} & \cdots & \frac{\partial y_{pq}}{\partial x_{ij}} \end{pmatrix}$

	$\frac{\partial \mathbf{Y}^{-1}}{\partial x} = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \mathbf{Y}^{-1}$		
--	---	--	--

4. Composite derivative

Given f and g are vector functions (or matrix functions), there are questions “how to take derivative of composite function $g \circ f = g(f(x))$ and how to extend basic derivative-taking techniques such as chain rule, product rule and sum rule from scalar function to vector function”. Let us sketch these rules before researching composite derivative in detailed.

Chain rule: Given the composite function $g \circ f = g(f(x))$ of two scalar functions f and g , the composite derivative is $(g \circ f)'(x) = g'(f(x)) = g'(y)f'(x) = \frac{\partial g \circ f}{\partial x} = \frac{\partial g}{\partial y} \frac{\partial y}{\partial x}$ where $y = f(x)$. Chain rule is applied to vector (or matrix) function in restriction. We make sense the compatibility between vector function f and g ; concretely, the image space of f must be the same to the domain space of g and the product of two derivatives g' and f' must be defined. For example, given vector

$$\text{variables } \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_p \end{pmatrix}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} \text{ and } \mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{pmatrix} \text{ where } \mathbf{z} \text{ is function of } \mathbf{y} \text{ which in turn is function of } \mathbf{x}.$$

By applying the chain rule, the derivative of \mathbf{z} w.r.t \mathbf{x} is:

$$\begin{aligned} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} &= \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial z_1}{\partial y_1} & \frac{\partial z_1}{\partial y_2} & \cdots & \frac{\partial z_1}{\partial y_n} \\ \frac{\partial z_2}{\partial y_1} & \frac{\partial z_2}{\partial y_2} & \cdots & \frac{\partial z_2}{\partial y_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_m}{\partial y_1} & \frac{\partial z_m}{\partial y_2} & \cdots & \frac{\partial z_m}{\partial y_n} \end{pmatrix} \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_p} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \frac{\partial y_n}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_p} \end{pmatrix} \\ &= \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_1}{\partial x_2} & \cdots & \frac{\partial z_1}{\partial x_p} \\ \frac{\partial z_2}{\partial x_1} & \frac{\partial z_2}{\partial x_2} & \cdots & \frac{\partial z_2}{\partial x_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_m}{\partial x_1} & \frac{\partial z_m}{\partial x_2} & \cdots & \frac{\partial z_m}{\partial x_p} \end{pmatrix} \text{ where } \frac{\partial z_i}{\partial x_j} = \sum_{k=1}^n \frac{\partial z_i}{\partial y_k} \frac{\partial y_k}{\partial x_j} \end{aligned}$$

Product rule: Given the product fg of two scalar functions f and g , the derivative of this product is $(f(x)g(x))' = f'(x)g(x) + f(x)g'(x) = \frac{\partial f}{\partial x}g(x) + f(x)\frac{\partial g}{\partial x}$. The product rule can be applied to any vector and matrix function but it makes sense to define what *product* is and how the order of product is because there are many kinds of product such as scalar product between two vectors, matrix multiplication. Moreover matrix multiplication is not commutative and so its order is very important. For example, given two $n \times 1$ column vectors \mathbf{u} and \mathbf{v} , both of them are functions of scalar

x ; we have $\mathbf{u} = \begin{pmatrix} u_1(x) \\ u_2(x) \\ \vdots \\ u_n(x) \end{pmatrix}$ and $\mathbf{v} = \begin{pmatrix} v_1(x) \\ v_2(x) \\ \vdots \\ v_n(x) \end{pmatrix}$. By applying product rule, the derivative of product $\mathbf{u}^T \mathbf{v}$

resulting out a scalar is different from the one of product $\mathbf{u} \mathbf{v}^T$ resulting out a matrix.

$$\begin{aligned}
\frac{\partial \mathbf{u}^T \mathbf{v}}{\partial x} &= \mathbf{u}^T \frac{\partial \mathbf{v}}{\partial x} + \mathbf{v}^T \frac{\partial \mathbf{u}}{\partial x} = (u_1(x), u_2(x), \dots, u_n(x)) \begin{pmatrix} \frac{\partial v_1}{\partial x} \\ \frac{\partial v_2}{\partial x} \\ \vdots \\ \frac{\partial v_n}{\partial x} \end{pmatrix} + (v_1(x), v_2(x), \dots, v_n(x)) \begin{pmatrix} \frac{\partial u_1}{\partial x} \\ \frac{\partial u_2}{\partial x} \\ \vdots \\ \frac{\partial u_n}{\partial x} \end{pmatrix} \\
&= \sum_{i=1}^n \left(u_i(x) \frac{\partial v_i}{\partial x} + v_i(x) \frac{\partial u_i}{\partial x} \right) \\
\frac{\partial \mathbf{u} \mathbf{v}^T}{\partial x} &= \mathbf{u} \frac{\partial \mathbf{v}^T}{\partial x} + \mathbf{v} \frac{\partial \mathbf{u}^T}{\partial x} = \begin{pmatrix} u_1(x) \\ u_2(x) \\ \vdots \\ u_n(x) \end{pmatrix} \left(\frac{\partial v_1}{\partial x}, \frac{\partial v_2}{\partial x}, \dots, \frac{\partial v_n}{\partial x} \right) + \begin{pmatrix} v_1(x) \\ v_2(x) \\ \vdots \\ v_n(x) \end{pmatrix} \left(\frac{\partial u_1}{\partial x}, \frac{\partial u_2}{\partial x}, \dots, \frac{\partial u_n}{\partial x} \right) \\
&= \begin{pmatrix} z_{11} & z_{12} & \dots & z_{1n} \\ z_{21} & z_{22} & \dots & z_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ z_{n1} & z_{n2} & \dots & z_{nn} \end{pmatrix} \text{ where } z_{ij} = u_i(x) \frac{\partial v_j}{\partial x} + v_j(x) \frac{\partial u_i}{\partial x}
\end{aligned}$$

Sum rule: Given the sum $f + g$ of two scalar functions f and g , the derivative of this sum is $(f + g)'(x) = f'(x) + g'(x)$. Sum rule is applied to any vector and matrix function with condition that f and g belong to the same kind. For example, given two $m \times n$ matrices:

$$U = \begin{pmatrix} u_{11}(x) & u_{12}(x) & \dots & u_{1n}(x) \\ u_{21}(x) & u_{22}(x) & \dots & u_{2n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ u_{m1}(x) & u_{m2}(x) & \dots & u_{mn}(x) \end{pmatrix} \text{ and } V = \begin{pmatrix} v_{11}(x) & v_{12}(x) & \dots & v_{1n}(x) \\ v_{21}(x) & v_{22}(x) & \dots & v_{2n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ v_{m1}(x) & v_{m2}(x) & \dots & v_{mn}(x) \end{pmatrix}$$

which are matrix functions of scalar x , the derivative of their sum is:

$$\begin{aligned}
\frac{\partial(\mathbf{U} + \mathbf{V})}{\partial x} &= \frac{\partial \mathbf{U}}{\partial x} + \frac{\partial \mathbf{V}}{\partial x} = \begin{pmatrix} \frac{\partial u_{11}}{\partial x} & \frac{\partial u_{12}}{\partial x} & \dots & \frac{\partial u_{1n}}{\partial x} \\ \frac{\partial u_{21}}{\partial x} & \frac{\partial u_{22}}{\partial x} & \dots & \frac{\partial u_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial u_{m1}}{\partial x} & \frac{\partial u_{m2}}{\partial x} & \dots & \frac{\partial u_{mn}}{\partial x} \end{pmatrix} + \begin{pmatrix} \frac{\partial v_{11}}{\partial x} & \frac{\partial v_{12}}{\partial x} & \dots & \frac{\partial v_{1n}}{\partial x} \\ \frac{\partial v_{21}}{\partial x} & \frac{\partial v_{22}}{\partial x} & \dots & \frac{\partial v_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial v_{m1}}{\partial x} & \frac{\partial v_{m2}}{\partial x} & \dots & \frac{\partial v_{mn}}{\partial x} \end{pmatrix} \\
&= \begin{pmatrix} \frac{\partial u_{11}}{\partial x} + \frac{\partial v_{11}}{\partial x} & \frac{\partial u_{12}}{\partial x} + \frac{\partial v_{12}}{\partial x} & \dots & \frac{\partial u_{1n}}{\partial x} + \frac{\partial v_{1n}}{\partial x} \\ \frac{\partial u_{21}}{\partial x} + \frac{\partial v_{21}}{\partial x} & \frac{\partial u_{22}}{\partial x} + \frac{\partial v_{22}}{\partial x} & \dots & \frac{\partial u_{2n}}{\partial x} + \frac{\partial v_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial u_{m1}}{\partial x} + \frac{\partial v_{m1}}{\partial x} & \frac{\partial u_{m2}}{\partial x} + \frac{\partial v_{m2}}{\partial x} & \dots & \frac{\partial u_{mn}}{\partial x} + \frac{\partial v_{mn}}{\partial x} \end{pmatrix}
\end{aligned}$$

By combination of these rules, it is easy for us to take derivative of any composite function if some conditions w.r.t vector (matrix) functions are satisfied. Please pay attention that numerator layout convention for derivative notation is followed. Moreover, given a function and one of its variables so-called x , the constant w.r.t variable x is defined as a scalar, vector or matrix whose evaluation is

not a function of x but such constant can be a function of other variables. There are three kinds of constant such as scalar constant, vector constant and matrix constant.

Next section will discuss composite derivative of vector and matrix function, hence, derivative of simple scalar function is ignored but its variants related to vector and matrix will be surveyed later.

Scalar-by-vector composite derivative

Following are composite derivatives of typical scalar-by-vector functions (Wikipedia, Matrix calculus, 2014). Note that independent variable x is always column vector and derivative of scalar-by-vector function results out gradient vector.

Condition	Scalar-by-vector derivative
a is a scalar constant	$\frac{\partial a}{\partial x} = \mathbf{0}^T$
a is a scalar constant $u = u(x)$ is a scalar function of x	$\frac{\partial au}{\partial x} = a \frac{\partial u}{\partial x}$
$u = u(x), v = v(x)$ are scalar functions of x	$\frac{\partial(u + v)}{\partial x} = \frac{\partial u}{\partial x} + \frac{\partial v}{\partial x}$
$u = u(x), v = v(x)$ are scalar functions of x	$\frac{\partial uv}{\partial x} = u \frac{\partial v}{\partial x} + v \frac{\partial u}{\partial x}$
$u = u(x)$ is a scalar function of x f is a scalar-by-scalar function	$\frac{\partial f(u)}{\partial x} = \frac{\partial f}{\partial u} \frac{\partial u}{\partial x}$
$u = u(x)$ is a scalar function of x f, g are scalar-by-scalar functions	$\frac{\partial g(f(u))}{\partial x} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial u} \frac{\partial u}{\partial x}$
a is a vector constant	$\frac{\partial a^T x}{\partial x} = \frac{\partial x^T a}{\partial x} = a^T$
	$\frac{\partial x^T x}{\partial x} = 2x^T$
$u = u(x), v = v(x)$ are vector functions of x	$\frac{\partial u^T v}{\partial x} = u^T \frac{\partial v}{\partial x} + v^T \frac{\partial u}{\partial x}$
	Note that $\frac{\partial v}{\partial x}$ and $\frac{\partial u}{\partial x}$ are derivatives of vector-by-scalar functions, which results out Jacobian matrix, which implies that $u^T \frac{\partial v}{\partial x}$ and $v^T \frac{\partial u}{\partial x}$ are vectors.

$\mathbf{u} = \mathbf{u}(x)$ is a vector function of x	$\frac{\partial a^T \mathbf{u}}{\partial x} = \frac{\partial \mathbf{u}^T a}{\partial x} = a^T \frac{\partial \mathbf{u}}{\partial x}$
a and b are vector constants	$\frac{\partial a^T x x^T b}{\partial x} = x^T (ab^T + ba^T)$
b is a vector constant. A is a matrix constant	$\frac{\partial b^T A x}{\partial x} = \frac{\partial x^T A^T b}{\partial x} = b^T A$
A is a square matrix constant	$\frac{\partial x^T A x}{\partial x} = x^T (A + A^T)$
A is a symmetric matrix constant	$\frac{\partial x^T A x}{\partial x} = 2x^T A$
A is a square matrix constant	$\frac{\partial^2 x^T A x}{\partial x^2} = A + A^T$
A is a symmetric matrix constant	$\frac{\partial^2 x^T A x}{\partial x^2} = 2A$
$\mathbf{u} = \mathbf{u}(x)$, $\mathbf{v} = \mathbf{v}(x)$ are vector functions of x and A is matrix constant.	$\frac{\partial \mathbf{u}^T A \mathbf{v}}{\partial x} = \mathbf{u}^T A \frac{\partial \mathbf{v}}{\partial x} + \mathbf{v}^T A^T \frac{\partial \mathbf{u}}{\partial x}$
A , C and D are matrix constants b and e are vector constants	$\begin{aligned} \frac{\partial (Ax + b)^T C(Dx + e)}{\partial x} &= \\ &\frac{\partial \mathbf{x}}{\partial x} (Dx + e)^T C^T A + (Ax + b)^T CD \end{aligned}$
a and a are vector constants	$\frac{\partial x - a }{\partial x} = \frac{(x - a)^T}{ x - a }$ Note that $ x - a $ is the length or norm or module of vector $x - a$

Vector-by-scalar composite derivative

Following are composite derivatives of typical vector-by-scalar functions (Wikipedia, Matrix calculus, 2014). Note that independent variable x is always scalar and derivative of scalar-by-vector function results out tangent vector.

Condition	Vector-by-scalar derivative
a is a vector constant	$\frac{\partial a}{\partial x} = \mathbf{0}$
a is a scalar constant $\mathbf{u} = \mathbf{u}(x)$ is a vector function of x	$\frac{\partial a \mathbf{u}}{\partial x} = a \frac{\partial \mathbf{u}}{\partial x}$
A is a matrix constant $\mathbf{u} = \mathbf{u}(x)$ is a vector function of x	$\frac{\partial A \mathbf{u}}{\partial x} = A \frac{\partial \mathbf{u}}{\partial x}$

$\mathbf{u} = \mathbf{u}(x)$ is a vector function of x	$\frac{\partial \mathbf{u}^T}{\partial x} = \left(\frac{\partial \mathbf{u}}{\partial x} \right)^T$
$\mathbf{u} = \mathbf{u}(x), \mathbf{v} = \mathbf{v}(x)$ are vector functions of x	$\frac{\partial(\mathbf{u} + \mathbf{v})}{\partial x} = \frac{\partial \mathbf{u}}{\partial x} + \frac{\partial \mathbf{v}}{\partial x}$
$\mathbf{u} = \mathbf{u}(x), \mathbf{v} = \mathbf{v}(x)$ are vector functions of x	$\frac{\partial(\mathbf{u} \times \mathbf{v})}{\partial x} = \mathbf{u} \times \frac{\partial \mathbf{v}}{\partial x} + \frac{\partial \mathbf{u}}{\partial x} \times \mathbf{v}$ <p>Note that sign \times represents cross product between two vectors, mentioned in previous section.</p>
$\mathbf{u} = \mathbf{u}(x)$ is a vector function of x f is a vector-by-vector function	$\frac{\partial f(\mathbf{u})}{\partial x} = \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$ <p>Note that vector-by-vector derivative gives Jacobian matrix, which refers that the product $\frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$ results out a vector.</p>
$\mathbf{u} = \mathbf{u}(x)$ is a scalar function of x f, g are vector-by-vector functions	$\frac{\partial g(f(\mathbf{u}))}{\partial x} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$

Vector-by-vector composite derivative

Following are composite derivatives of typical vector-by-vector functions (Wikipedia, Matrix calculus, 2014). Note that and independent variable x is always column vector and derivative of vector-by-vector function results out Jacobian matrix.

Condition	Vector-by-vector derivative
a is a vector constant	$\frac{\partial a}{\partial x} = \mathbf{0}$
	$\frac{\partial x}{\partial x} = I$
A is a matrix constant	$\frac{\partial Ax}{\partial x} = A$
A is a matrix constant	$\frac{\partial x^T A}{\partial x} = A^T$
a is a scalar constant $\mathbf{u} = \mathbf{u}(x)$ is a vector function of x	$\frac{\partial au}{\partial x} = a \frac{\partial u}{\partial x}$
$a = a(x)$ is a scalar function of x $\mathbf{u} = \mathbf{u}(x)$ is a vector function of x	$\frac{\partial au}{\partial x} = a \frac{\partial u}{\partial x} + u \frac{\partial a}{\partial x}$

	Note that scalar-by-scalar derivative $\frac{\partial a}{\partial x}$ gives out gradient vector (row vector) and so the expression $a \frac{\partial u}{\partial x} + u \frac{\partial a}{\partial x}$ results out a matrix
A is matrix constant $\mathbf{u} = \mathbf{u}(x)$ is a vector function of x	$\frac{\partial A\mathbf{u}}{\partial x} = A \frac{\partial \mathbf{u}}{\partial x}$
$\mathbf{u} = \mathbf{u}(x), \mathbf{v} = \mathbf{v}(x)$ are vector functions of x	$\frac{\partial(\mathbf{u} + \mathbf{v})}{\partial x} = \frac{\partial \mathbf{u}}{\partial x} + \frac{\partial \mathbf{v}}{\partial x}$
$\mathbf{u} = \mathbf{u}(x)$ is a vector function of x f is a vector-by-vector function	$\frac{\partial f(\mathbf{u})}{\partial x} = \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$ Note that vector-by-vector derivative gives Jacobian matrix, which refers that the matrix multiplication $\frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$ results out a matrix.
$\mathbf{u} = \mathbf{u}(x)$ is a vector function of x f, g are vector-by-vector functions	$\frac{\partial g(f(\mathbf{u}))}{\partial x} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$

Matrix function and differential

Before researching matrix derivative, we should have a preparation of concepts of matrix function and differential. Among scalar-by-matrix functions, determinant $|X|$ and operator trace $tr(X)$ are the most important. Given square matrix X , the derivative of determinant $|X|$ is adjoint matrix $adj(X)$; please see section “Basic concepts” for adjoint matrix definition.

$$\frac{\partial |X|}{\partial X} = adj(X) = |X|X^{-1}$$

Given square matrix, the derivative of operator trace $tr(X)$ is identity matrix I .

$$\frac{\partial tr(X)}{\partial X} = I$$

Differential is independent concept but it is often defined via derivative. Given matrix variable X , differential of X denoted dX is an infinitesimal according to pre-defined metric. If X is defined in real number field \mathbb{R} , then differential dX is

$$dX = \begin{pmatrix} dx_{11} & dx_{12} & \dots & dx_{1n} \\ dx_{21} & dx_{22} & \dots & dx_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ dx_{m1} & dx_{m2} & \dots & dx_{mn} \end{pmatrix} \text{ where } dx_{ij} \text{ is differential of real number } x_{ij}$$

Given matrix function $Y = F(X)$, differential of Y denoted dY is defined as multiplication of derivative of Y w.r.t X by differential of X .

$$d\mathbf{Y} = d\mathbf{F} = \mathbf{F}'(\mathbf{X})d\mathbf{X} = \frac{\partial \mathbf{Y}}{\partial \mathbf{X}}d\mathbf{X}$$

For convenience, we use derivative notation for denoting differential. It means that differentials of \mathbf{X} and \mathbf{Y} are denoted $\partial\mathbf{X}$ and $\partial\mathbf{Y}$, respectively instead of $d\mathbf{X}$ and $d\mathbf{Y}$. Note that vector is reduced form of matrix and so, differential of vector is defined by the same way. For example, given vector-by-vector function $\mathbf{y} = \mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ where \mathbf{A} is matrix constant, derivative of vector \mathbf{y} w.r.t vector \mathbf{x} is:

$$\frac{\partial \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}$$

The differential of \mathbf{y} is:

$$\partial\mathbf{y} = \mathbf{A}\partial\mathbf{x}$$

Therefore, the concept of differential totally coincides with the concept of derivative. The semantic of differential is deviation and differential is used to approximate a function or dependent variable. The deviation of \mathbf{y} denoted $\partial\mathbf{y}$ is approximated by the deviation of \mathbf{x} denoted $\partial\mathbf{x}$ and so the approximation expression is $\partial\mathbf{y} = \mathbf{A}\partial\mathbf{x}$ in which \mathbf{A} is result of derivative of \mathbf{y} . Two expressions $\frac{\partial \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}$ and $\partial\mathbf{y} = \mathbf{A}\partial\mathbf{x}$ are the same. As a result, all aforementioned derivatives can be written in differential form. If you want to convert a complicated differential to a derivative and vice versa, please find out the canonical form of such differential and look up the respective derivative of canonical form in following table (Wikipedia, Matrix calculus, 2014).

Canonical differential form	Equivalent derivative form	Note
$\partial\mathbf{y} = a\partial\mathbf{x}$	$\frac{\partial\mathbf{y}}{\partial\mathbf{x}} = a$	y is dependent scalar variable. x is independent scalar variable. a is scalar.
$\partial\mathbf{y} = \mathbf{a}\partial\mathbf{x}$	$\frac{\partial\mathbf{y}}{\partial\mathbf{x}} = \mathbf{a}$	y is dependent scalar variable. \mathbf{x} is independent vector variable. \mathbf{a}^T is gradient (row) vector.
$\partial\mathbf{y} = \text{tr}(\mathbf{A}\partial\mathbf{X})$	$\frac{\partial\mathbf{y}}{\partial\mathbf{X}} = \mathbf{A}$	y is dependent scalar variable. \mathbf{X} is independent matrix variable. \mathbf{A} is gradient matrix.
$\partial\mathbf{y} = \mathbf{a}\partial\mathbf{x}$	$\frac{\partial\mathbf{y}}{\partial\mathbf{x}} = \mathbf{a}$	y is dependent vector variable. x is independent scalar variable. \mathbf{a} is tangent (column) vector.
$\partial\mathbf{y} = \mathbf{A}\partial\mathbf{x}$	$\frac{\partial\mathbf{y}}{\partial\mathbf{x}} = \mathbf{A}$	y is dependent vector variable. \mathbf{x} is independent vector variable. \mathbf{A} is Jacobian matrix.
$\partial\mathbf{Y} = \mathbf{A}\partial\mathbf{x}$	$\frac{\partial\mathbf{Y}}{\partial\mathbf{x}} = \mathbf{A}$	\mathbf{Y} is dependent matrix variable. x is independent scalar variable. \mathbf{A} is tangent matrix.

All of principles and techniques in matrix derivative such as chain rule, product rule, sum rule, and etc are applied in matrix differential. Let \mathbf{X} and \mathbf{Y} be dependent square matrix variables, following are some basic properties of differential inferred directly from derivatives (Petersen & Pedersen, 2012, p. 8).

$$\partial A = 0 \text{ where } A \text{ is constant}$$

$$\partial \mathbf{X}^T = (\partial \mathbf{X})^T$$

$$\partial a\mathbf{X} = a\partial \mathbf{X} \text{ where } a \text{ is constant}$$

$$\partial \text{tr}(\mathbf{X}) = \text{tr}(\partial \mathbf{X})$$

$$\partial(\mathbf{X} + \mathbf{Y}) = \partial \mathbf{X} + \partial \mathbf{Y}$$

$$\partial \mathbf{XY} = (\partial \mathbf{X})\mathbf{Y} + \mathbf{X}(\partial \mathbf{Y})$$

$$\partial \mathbf{X}^{-1} = -\mathbf{X}^{-1}\partial \mathbf{XX}^{-1}$$

$$\partial |\mathbf{X}| = \text{tr}(\text{adj}(\mathbf{X})\partial \mathbf{X}) = |\mathbf{X}| \text{tr}(\mathbf{X}^{-1}\partial \mathbf{X})$$

where $\text{adj}(\mathbf{X})$ is adjoint matrix of \mathbf{X}

$$\partial(\ln|\mathbf{X}|) = \text{tr}(\mathbf{X}^{-1}\partial \mathbf{X})$$

We have an example of trace differential for illustrating the mutual relationship of differential and derivative. Let $\mathbf{Y} = \mathbf{AXBX}^T\mathbf{C}$ be matrix function by matrix variable \mathbf{X} , all of them are square matrices. The differential of \mathbf{Y} (Wikipedia, Matrix calculus, 2014) is:

$$\begin{aligned} \partial \text{tr}(\mathbf{AXBX}^T\mathbf{C}) &= \partial \text{tr}(\mathbf{CAXBX}^T) = \text{tr}(\partial(\mathbf{CAXBX}^T)) = \text{tr}(\partial(\mathbf{CAX})\mathbf{BX}^T + \mathbf{CAX}\partial(\mathbf{BX}^T)) \\ &= \text{tr}(\mathbf{CA}\partial(\mathbf{X})\mathbf{BX}^T) + \text{tr}(\mathbf{CAXB}\partial \mathbf{X}^T) = \text{tr}(\mathbf{BX}^T\mathbf{CA}\partial(\mathbf{X})) + \text{tr}(\mathbf{CAXB}(\partial \mathbf{X})^T) \\ &= \text{tr}(\mathbf{BX}^T\mathbf{CA}\partial(\mathbf{X})) + \text{tr}((\mathbf{CAXB}(\partial \mathbf{X})^T)^T) \\ &= \text{tr}(\mathbf{BX}^T\mathbf{CA}\partial(\mathbf{X})) + \text{tr}(\partial \mathbf{XB}^T\mathbf{X}^T\mathbf{A}^T\mathbf{C}^T) = \text{tr}(\mathbf{BX}^T\mathbf{CA}\partial(\mathbf{X})) + \text{tr}(\mathbf{B}^T\mathbf{X}^T\mathbf{A}^T\mathbf{C}^T\partial \mathbf{X}) \\ &= \text{tr}((\mathbf{BX}^T\mathbf{CA} + \mathbf{B}^T\mathbf{X}^T\mathbf{A}^T\mathbf{C}^T)\partial \mathbf{X}) \end{aligned}$$

$$\Rightarrow \frac{\partial \text{tr}(\mathbf{AXBX}^T\mathbf{C})}{\partial \mathbf{X}} = \mathbf{BX}^T\mathbf{CA} + \mathbf{B}^T\mathbf{X}^T\mathbf{A}^T\mathbf{C}^T$$

We will survey matrix-by-matrix functions. Aforementioned in section “Matrix analysis”, given non-singular matrix $\mathbf{A}(n,n)$, its eigenvalues are solutions of following equation:

$$|\mathbf{A} - \lambda \mathbf{I}_n| = 0 \text{ where } \mathbf{I}_n \text{ is } n,n \text{ identity matrix}$$

Determinant $|\mathbf{A} - \lambda \mathbf{I}_n|$ is expanded as n^{th} order polynomial called *characteristic polynomial* in which λ is scalar variable, we re-write this polynomial.

$$P(\lambda) = |\mathbf{A} - \lambda \mathbf{I}_n|$$

Cayley–Hamilton theorem (Wikipedia, Cayley–Hamilton theorem, 2014) states that if $P(\lambda)$ has solution(s) then, when substituting λ by matrix \mathbf{A} , we get a zero matrix $\mathbf{0}$.

$$P(\mathbf{A}) = \mathbf{0}$$

For example let $\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$, the characteristic polynomial is expressed as below:

$$P(\lambda) = |A - \lambda I_2| = \begin{vmatrix} 1-\lambda & 2 \\ 3 & 4-\lambda \end{vmatrix} = \lambda^2 - 5\lambda - 2$$

Because $P(\lambda) = 0$ has solutions $\lambda = \frac{5 \pm \sqrt{33}}{2}$, we have:

$$\mathbf{P}(A) = A^2 - 5A - 2I_2 = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} - 5 \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} - 2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$$

Given any square matrix constant C , matrix polynomial $2C = \mathbf{P}(A) = A^2 - 5A - 2C$ and its derivative always exist.

$$\frac{\partial \mathbf{P}}{\partial A} = 2A - 5I_2$$

In generalization, given matrix variable X , the **matrix polynomial** known as matrix-by-matrix function always exists together its derivative if X is non-singular square matrix. Note that matrix operations within matrix polynomial include scalar multiplication, addition, subtraction, multiplication and power.

$$\mathbf{P}(X) = a_0 I + a_1 X + a_2 X^2 + \cdots + a_n X^n = \sum_{k=0}^n a_k X^k$$

$$\mathbf{P}'(X) = \frac{\partial \mathbf{P}}{\partial X} = a_1 I + 2a_2 X + 3a_3 X^2 + \cdots + n a_n X^{n-1} = \sum_{k=1}^n k a_k X^{k-1}$$

Where a_i (s) are scalar and I is identity matrix.

Given real scalar function f has Taylor expansion or Taylor series at entry 0.

$$f(x) = f(0) + \frac{f'(0)x}{1!} + \frac{f''(0)x^2}{2!} + \cdots + \frac{f^{(n)}(0)x^n}{n!} + \cdots = \sum_{k=0}^{\infty} \frac{f^{(k)}(0)x^k}{k!}$$

where $f^{(k)}$ denotes k^{th} derivative of f and $f^{(0)} = f$

When substituting x by square matrix X , we get an matrix-by-matrix function called **matrix-power-series function** in which arithmetic operations are extended as matrix operations such as scalar multiplication, addition, subtraction, multiplication and power. Following is Taylor series of matrix-power-series function S at matrix entry 0.

$$S(X) = S(0) + \frac{S'(0)X}{1!} + \frac{S''(0)X^2}{2!} + \cdots + \frac{S^{(n)}(0)X^n}{n!} + \cdots = \sum_{k=0}^{\infty} \frac{S^{(k)}(0)X^k}{k!}$$

where $S^{(k)}$ denotes k^{th} derivative of S and $S^{(0)} = S$

There are some popular matrix-power-series functions such as matrix exponential E^X , natural logarithm of matrix LNX , $SINX$, $COSX$. Because all matrix-power-series functions $S(X)$ are

defined by Taylor series, it is very complicated to evaluate them and so there is a practical method

to calculate them. If $X = \begin{pmatrix} x_{11} & 0 & \dots & 0 \\ 0 & x_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & x_{nn} \end{pmatrix}$ is diagonal matrix then

$$S(X) = \begin{pmatrix} s(x_{11}) & 0 & \dots & 0 \\ 0 & s(x_{22}) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s(x_{nn}) \end{pmatrix}$$

Where $S(X)$ represents matrix-power-series functions such as matrix exponential E^X , natural logarithm of matrix $LN X$, $SIN X$, $COS X$; s represents scalar functions such as e , \ln , \sin , \cos . If diagonalizable matrix X is decomposed into spectrums $X = UAU^{-1}$ then

$$S(X) = US(A)U^{-1}$$

The formula above is very important for calculating matrix-power-series function; which implies it is possible to compute practically any diagonalizable matrix instead of using complicated Taylor series. For example, given 2x2 matrix $X = \begin{pmatrix} 3/2 & 1/2 \\ 1/2 & 3/2 \end{pmatrix}$ is diagonalized as follows:

$$X = \begin{pmatrix} 3/2 & 1/2 \\ 1/2 & 3/2 \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix}$$

Matrix-power-series functions E^X , $LN X$, $SIN X$ and $COS X$ are totally evaluated.

$$\begin{aligned} E^X &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} E \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \\ &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \begin{pmatrix} e^1 & 0 \\ 0 & e^2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} = \frac{e}{2} \begin{pmatrix} e+1 & e-1 \\ e-1 & e+1 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} LN X &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} LN \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \\ &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \begin{pmatrix} \ln 1 & 0 \\ 0 & \ln 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} = \frac{\ln 2}{2} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} SIN X &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} SIN \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \\ &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \begin{pmatrix} \sin 1 & 0 \\ 0 & \sin 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} \sin 2 + \sin 1 & \sin 2 - \sin 1 \\ \sin 2 - \sin 1 & \sin 2 + \sin 1 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} COS X &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} COS \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \\ &= \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \begin{pmatrix} \cos 1 & 0 \\ 0 & \cos 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} \cos 2 + \cos 1 & \cos 2 - \cos 1 \\ \cos 2 - \cos 1 & \cos 2 + \cos 1 \end{pmatrix} \end{aligned}$$

Now we research popular matrix-power-series functions such as matrix exponential \mathbf{E}^X , natural logarithm of matrix $\mathbf{LN}X$, $\mathbf{SIN}X$ and $\mathbf{COS}X$. The exponential of square matrix X denoted \mathbf{E}^X or $\mathbf{EXP}(X)$ is invertible matrix which defined as Taylor series (Wikipedia, Matrix exponential, 2014):

$$\mathbf{E}^X = \mathbf{I} + \frac{\mathbf{X}}{1!} + \frac{\mathbf{X}^2}{2!} + \cdots + \frac{\mathbf{X}^n}{n!} + \cdots = \sum_{k=0}^{\infty} \frac{\mathbf{X}^k}{k!}$$

and its derivative is

$$(\mathbf{E}^X)' = \frac{\partial \mathbf{E}^X}{\partial \mathbf{X}} = \mathbf{E}^X$$

Given $x \in \mathbb{R}$ is real scalar and A is matrix constant, we have

$$\frac{\partial \mathbf{E}^{xA}}{\partial x} = A \mathbf{E}^{tA} = \mathbf{E}^{tA} A$$

Given square matrix $\mathbf{U}(x) = \begin{pmatrix} u_{11}(x) & u_{12}(x) & \dots & u_{1n}(x) \\ u_{21}(x) & u_{22}(x) & \dots & u_{2n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ u_{m1}(x) & u_{m2}(x) & \dots & u_{mn}(x) \end{pmatrix}$ is matrix-by-scalar function, the derivative of $\mathbf{E}^{\mathbf{U}(x)}$ w.r.t scalar variable x is (Wikipedia, Matrix exponential, 2014):

$$\frac{\partial \mathbf{E}^{\mathbf{U}(x)}}{\partial x} = \int_0^1 \mathbf{E}^{t\mathbf{U}(x)} \frac{\partial \mathbf{U}}{\partial x} \mathbf{E}^{(1-t)\mathbf{U}(x)} dt$$

Following are other properties of \mathbf{E}^X (Wikipedia, Matrix exponential, 2014):

\mathbf{E}^X is always invertible matrix

$$\mathbf{E}^0 = \mathbf{I}$$

$$\mathbf{E}^{aX} \mathbf{E}^{bX} = \mathbf{E}^{(a+b)X}$$

$$\mathbf{E}^X \mathbf{E}^{-X} = \mathbf{I}$$

If X and Y are mutually commutative $XY = YX$ then $\mathbf{E}^X \mathbf{E}^Y = \mathbf{E}^Y \mathbf{E}^X = \mathbf{E}^{(X+Y)}$

If Y is invertible then $\mathbf{E}^{YXY^{-1}} = \mathbf{Y} \mathbf{E}^X \mathbf{Y}^{-1}$

$$\mathbf{EXP}(X^T) = (\mathbf{EXP}(X))^T$$

$$|\mathbf{E}^X| = \mathbf{E}^{tr(X)}$$

If X is idempotent matrix ($XX = X$) then $\mathbf{E}^X = \mathbf{I} + (e - 1)X$

The natural logarithm (Wikipedia, Logarithm of a matrix, 2014) of invertible matrix X denoted $\mathbf{LN}X$ is defined as the matrix Y such that $Y = \mathbf{E}^X$. Note that $\mathbf{LN}X$ exists if and only if X is invertible matrix. $\mathbf{LN}X$ is defined as Taylor series:

$$\mathbf{LN}(X) = \mathbf{LN}(\mathbf{I} + A) = A - \frac{\mathbf{A}^2}{2} + \frac{\mathbf{A}^3}{3} - \frac{\mathbf{A}^4}{4} + \cdots = \sum_{k=0}^{\infty} \frac{(-1)^k \mathbf{A}^{k+1}}{k+1}$$

Following are properties of $\mathbf{LN}X$ (Wikipedia, Logarithm of a matrix, 2014):

$$X^{-1} = e^{-\ln(X)}$$

If X and Y are both positive-definite and commutative matrices, then

$$XY = \mathbf{E}^{\mathbf{LN}X + \mathbf{LN}Y}$$

Given square matrix X , sine and cosine of X is defined by Taylor series.

$$\sin X = X - \frac{X^3}{3!} + \frac{X^5}{5!} - \frac{X^7}{7!} + \dots + \frac{(-1)^n X^{2n+1}}{(2n+1)!} + \dots = \sum_{k=0}^{\infty} \frac{(-1)^k X^{2k+1}}{(2k+1)!}$$

$$\cos X = I - \frac{X^2}{2!} + \frac{X^4}{4!} - \frac{X^6}{6!} + \dots + \frac{(-1)^n X^{2n}}{(2n)!} + \dots = \sum_{k=0}^{\infty} \frac{(-1)^k X^{2k}}{(2k)!}$$

The derivatives of sine and cosine of X are:

$$\begin{aligned}\sin' X &= \frac{\partial \sin X}{\partial X} = \cos X \\ \cos' X &= \frac{\partial \cos X}{\partial X} = -\sin X\end{aligned}$$

Following table is sketchy summarization of matrix functions, suppose that independent variable X is square matrix.

Dependent variable	Independent variable		
	Scalar	Vector	Matrix
Scalar			Determinant $ X $ with X is square matrix $\frac{\partial X }{\partial X} = adj(X) = X X^{-1}$
			Trace operator $tr(X)$ $\frac{\partial tr(X)}{\partial X} = I$ $\partial tr(X) = tr(\partial X)$ $\partial X = tr(adj(X)\partial X) = X tr(X^{-1}\partial X)$
Vector			Matrix polynomial $P(X) = \sum_{k=0}^n a_k X^k$ $P'(X) = \frac{\partial P}{\partial X} = \sum_{k=1}^n k a_k X^{k-1}$
Matrix			Matrix-power-series function such as E^X , $\ln X$, $\sin X$ and $\cos X$ $E^X = \sum_{k=0}^{\infty} \frac{X^k}{k!}$ $(E^X)' = \frac{\partial E^X}{\partial X} = E^X$ $\sin X = \sum_{k=0}^{\infty} \frac{(-1)^k X^{2k+1}}{(2k+1)!}$ $\cos X = \sum_{k=0}^{\infty} \frac{(-1)^k X^{2k}}{(2k)!}$ $\sin' X = \cos X$

			$\cos'X = -\sin X$
--	--	--	--------------------------------------

Other composite derivatives relevant to matrix will be described below.

Scalar-by-matrix composite derivative

Determinant and trace operators are important to scalar-by-matrix functions and hence, scalar-by-matrix derivatives often relate to determinant and trace operators. Following are composite derivatives of typical scalar-by-matrix functions. Note that independent variable X is always matrix and derivative of scalar-by-matrix function results out gradient matrix.

Condition	Scalar-by-matrix derivative
a is a scalar constant	$\frac{\partial a}{\partial X} = \mathbf{0}^T$
a is a scalar constant $u = u(X)$ is a scalar function of X	$\frac{\partial au}{\partial X} = a \frac{\partial u}{\partial X}$
$u = u(X), v = v(X)$ are scalar functions of X	$\frac{\partial(u + v)}{\partial X} = \frac{\partial u}{\partial X} + \frac{\partial v}{\partial X}$
$u = u(X), v = v(X)$ are scalar functions of X	$\frac{\partial uv}{\partial X} = u \frac{\partial v}{\partial X} + v \frac{\partial u}{\partial X}$
$u = u(X)$ is a scalar function of X f is a scalar-by-scalar function	$\frac{\partial f(u)}{\partial X} = \frac{\partial f}{\partial u} \frac{\partial u}{\partial X}$
$u = u(X)$ is a scalar function of X f, g are scalar-by-scalar functions	$\frac{\partial g(f(u))}{\partial X} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial u} \frac{\partial u}{\partial X}$
$U = U(X)$ is a matrix function of X f is scalar-by-matrix function	$\frac{\partial f(U)}{\partial x_{ij}} = \text{tr} \left(\frac{\partial f(U)}{\partial U} \frac{\partial U}{\partial x_{ij}} \right)$ Note that $\frac{\partial f(U)}{\partial x_{ij}}$ is derivative of scalar-by-scalar function, which results out a scalar

Trace related derivative

	$\frac{\partial \text{tr}(X)}{\partial X} = I$
$G(X)$ is any matrix polynomial $P(X)$ or matrix-power-series function such as E^X ,	Note that the derivative of trace operator is mentioned in previous section $\frac{\partial \text{tr}(G(X))}{\partial X} = G'(X)$

LNX, $SINX$, $COSX$ and etc	<p>Following are some popular derivatives</p> $\frac{\partial \text{tr}(\mathbf{P}(\mathbf{X}))}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\sum_{k=0}^n a_k \mathbf{X}^k)}{\partial \mathbf{X}} = \mathbf{P}'(\mathbf{X})$ $= \sum_{k=1}^n k a_k \mathbf{X}^{k-1}$ $\frac{\partial \text{tr}(\mathbf{X}^n)}{\partial \mathbf{X}} = n \mathbf{X}^{n-1} \text{ where } n \text{ is positive integer}$ $\frac{\partial \text{tr}(\mathbf{E}^{\mathbf{X}})}{\partial \mathbf{X}} = \mathbf{E}^{\mathbf{X}}$ $\frac{\partial \text{tr}(\mathbf{SINX})}{\partial \mathbf{X}} = \mathbf{COSX}$ $\frac{\partial \text{tr}(\mathbf{COSX})}{\partial \mathbf{X}} = -\mathbf{SINX}$
$\mathbf{U} = \mathbf{U}(\mathbf{X})$ and $\mathbf{V} = \mathbf{V}(\mathbf{X})$ are matrix functions of \mathbf{X}	$\frac{\partial \text{tr}(\mathbf{U} + \mathbf{V})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{U})}{\partial \mathbf{X}} + \frac{\partial \text{tr}(\mathbf{V})}{\partial \mathbf{X}}$
a is a scalar constant $\mathbf{U} = \mathbf{U}(\mathbf{X})$ is a matrix function of \mathbf{X}	$\frac{\partial \text{tr}(a\mathbf{U})}{\partial \mathbf{X}} = a \frac{\partial \text{tr}(\mathbf{U})}{\partial \mathbf{X}}$
A, B are matrix constants	$\frac{\partial \text{tr}(\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{X}\mathbf{A})}{\partial \mathbf{X}} = A$ $\frac{\partial \text{tr}(\mathbf{A}\mathbf{X}\mathbf{B})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{B}\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = BA$ $\frac{\partial \text{tr}(\mathbf{A}\mathbf{X}^T)}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{X}^T\mathbf{A})}{\partial \mathbf{X}} = A^T$ $\frac{\partial \text{tr}(\mathbf{A}\mathbf{X}^T\mathbf{B})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{B}\mathbf{A}\mathbf{X}^T)}{\partial \mathbf{X}} = A^T B^T$
A is a matrix constant	$\frac{\partial \text{tr}(\mathbf{X}^T\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{X}\mathbf{A}\mathbf{X}^T)}{\partial \mathbf{X}} = \mathbf{X}^T(A + A^T)$
A, B are matrix constants	$\frac{\partial \text{tr}(\mathbf{X}^{-1}\mathbf{B})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{B}\mathbf{X}^{-1})}{\partial \mathbf{X}} = -(\mathbf{X}^{-1})^T B (\mathbf{X}^{-1})^T$ $\frac{\partial \text{tr}(\mathbf{A}\mathbf{X}^{-1}\mathbf{B})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{B}\mathbf{A}\mathbf{X}^{-1})}{\partial \mathbf{X}} = -(\mathbf{X}^{-1})^T B A (\mathbf{X}^{-1})^T$
A, B, C are matrix constants	$\frac{\partial \text{tr}(\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}^T\mathbf{C})}{\partial \mathbf{X}} = \mathbf{B}\mathbf{X}^T\mathbf{C}A + \mathbf{B}^T\mathbf{X}^TA^T\mathbf{C}^T$

A is matrix constant	$\frac{\partial \text{tr}(AX^n)}{\partial X} = \sum_{k=0}^{n-1} X^k AX^{n-k-1}$
<i>Determinant related derivative</i>	
	$\frac{\partial \mathbf{X} }{\partial \mathbf{X}} = \text{adj}(\mathbf{X}) = \mathbf{X} \mathbf{X}^{-1}$
a is scalar constant \mathbf{X}^- is G-inverse of \mathbf{X} , please see section “Basic concepts”	$\frac{\partial \ln a\mathbf{X} }{\partial \mathbf{X}} = \mathbf{X}^{-1}$ $\frac{\partial \ln \mathbf{X}^T \mathbf{X} }{\partial \mathbf{X}} = 2\mathbf{X}^-$ $\frac{\partial \ln \mathbf{X}^T \mathbf{X} }{\partial \mathbf{X}^-} = -2\mathbf{X}$
A, B are matrix constants	$\frac{\partial AXB }{\partial \mathbf{X}} = AXB \mathbf{X}^{-1}$
n is positive integer	$\frac{\partial \mathbf{X}^n }{\partial \mathbf{X}} = n \mathbf{X}^n \mathbf{X}^{-1}$
A is a matrix constant \mathbf{X} is invertible	$\frac{\partial \mathbf{X}^T A \mathbf{X} }{\partial \mathbf{X}} = 2 \mathbf{X}^T A \mathbf{X} \mathbf{X}^{-1}$
A is a matrix constant and symmetric \mathbf{X} is not square matrix	$\frac{\partial \mathbf{X}^T A \mathbf{X} }{\partial \mathbf{X}} = 2 \mathbf{X}^T A \mathbf{X} (\mathbf{X}^T A^T \mathbf{X})^{-1} \mathbf{X}^T A^T$
A is a matrix constant and non-symmetric \mathbf{X} is not square matrix	$\begin{aligned} \frac{\partial \mathbf{X}^T A \mathbf{X} }{\partial \mathbf{X}} &= 2 \mathbf{X}^T A \mathbf{X} ((\mathbf{X}^T A \mathbf{X})^{-1} \mathbf{X}^T A \\ &\quad + (\mathbf{X}^T A^T \mathbf{X})^{-1} \mathbf{X}^T A^T) \end{aligned}$

Matrix-by-scalar composite derivative

Following are composite derivatives of typical scalar-by-matrix functions. Note that independent variable x is always scalar and derivative of matrix-by-scalar function results out tangent matrix.

Condition	Matrix-by-scalar derivative
A is a matrix constant	$\frac{\partial A}{\partial x} = \mathbf{0}$
a is a scalar constant $\mathbf{U} = \mathbf{U}(x)$ is a matrix function of x	$\frac{\partial a\mathbf{U}}{\partial x} = a \frac{\partial \mathbf{U}}{\partial x}$

A, B are matrix constants $\mathbf{U} = \mathbf{U}(x)$ is a matrix function of x	$\frac{\partial A\mathbf{U}B}{\partial x} = A \frac{\partial \mathbf{U}}{\partial x} B$
$\mathbf{U} = \mathbf{U}(x), \mathbf{V} = \mathbf{V}(x)$ are matrix functions of x	$\frac{\partial (\mathbf{U} + \mathbf{V})}{\partial x} = \frac{\partial \mathbf{U}}{\partial x} + \frac{\partial \mathbf{V}}{\partial x}$
$\mathbf{U} = \mathbf{U}(x), \mathbf{V} = \mathbf{V}(x)$ are matrix functions of x	$\frac{\partial (\mathbf{U}\mathbf{V})}{\partial x} = \mathbf{U} \frac{\partial \mathbf{V}}{\partial x} + \frac{\partial \mathbf{U}}{\partial x} \mathbf{V}$
$\mathbf{U} = \mathbf{U}(x)$ is a matrix function of x	$\frac{\partial \mathbf{U}^{-1}}{\partial x} = -\mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \mathbf{U}^{-1}$
$\mathbf{U} = \mathbf{U}(x, y)$ is a matrix function of x and y . Note that \mathbf{U} is not scalar-by-matrix function, it is really vector-by-matrix function	$\frac{\partial \mathbf{U}^{-1}}{\partial x \partial y} = -\mathbf{U}^{-1} \left(\frac{\partial \mathbf{U}}{\partial x} \mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial y} - \frac{\partial^2 \mathbf{U}}{\partial x \partial y} \right. \\ \left. + \frac{\partial \mathbf{U}}{\partial y} \mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \right) \mathbf{U}^{-1}$
A is square matrix constant. $\mathbf{G}(\cdot)$ is any matrix polynomial $\mathbf{P}(X)$ or matrix-power-series function such as \mathbf{E}^X , $\mathbf{LN}X$, $\mathbf{SIN}X$, $\mathbf{COS}X$.	$\frac{\partial \mathbf{G}(xA)}{\partial x} = A\mathbf{G}'(xA) = \mathbf{G}'(xA)A$ Following are some popular derivatives $\frac{\partial \mathbf{P}(xA)}{\partial x} = \frac{\partial \sum_{k=0}^n a_k(xA)^k}{\partial x} = A\mathbf{P}'(xA)$ $= A \sum_{k=1}^n k a_k(xA)^{k-1}$ $= \left(\sum_{k=1}^n k a_k(xA)^{k-1} \right) A$ $\frac{\partial (xA)^n}{\partial x} = nA(xA)^{n-1}$ $= n(xA)^{n-1}A \text{ where } n \text{ is positive integer}$ $\frac{\partial \mathbf{E}^{xA}}{\partial x} = A\mathbf{E}^{xA} = \mathbf{E}^{xA}A$ $\frac{\partial \mathbf{SIN}(xA)}{\partial x} = A\mathbf{COS}(xA) = \mathbf{COS}(xA)A$ $\frac{\partial \mathbf{COS}(xA)}{\partial x} = -A\mathbf{SIN}(xA) = -\mathbf{SIN}(xA)A$
$\mathbf{U} = \mathbf{U}(x)$ is a matrix function of x and \mathbf{U} is square matrix. \mathbf{E}^X is matrix exponential function.	$\frac{\partial \mathbf{E}^{U(x)}}{\partial x} = \int_0^1 \mathbf{E}^{tU(x)} \frac{\partial \mathbf{U}}{\partial x} \mathbf{E}^{(1-t)U(x)} dt$

Scalar-by-scalar composite derivative with vector and matrix involved

There are many cases in which a scalar-by-scalar function is composed of vector and matrix functions, which requires some efforts to solve some new problems although scalar-by-scalar derivative is much simpler than vector and matrix derivative. Following are composite derivatives of typical scalar-by-scalar functions involved vector and matrix. Note that independent variable x is always scalar and derivative of scalar-by-scalar function results out scalar.

Condition	Vector-by-vector derivative
$\mathbf{u} = \mathbf{u}(x)$ is a vector-by-scalar function. f is scalar-by-vector function	$\frac{\partial f(\mathbf{u})}{\partial x} = \frac{\partial f(\mathbf{u})}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$ <p>Note that $\frac{\partial f(\mathbf{u})}{\partial \mathbf{u}}$ is row gradient vector and $\frac{\partial \mathbf{u}}{\partial x}$ is tangent vector. Therefore, the result $\frac{\partial f(\mathbf{u})}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial x}$ is really a scalar product of two vectors, which produces a scalar.</p>
$\mathbf{u} = \mathbf{u}(x), \mathbf{v} = \mathbf{v}(x)$ are vector-by-scalar functions	$\frac{\partial \mathbf{u}^T \mathbf{v}}{\partial x} = \mathbf{u}^T \frac{\partial \mathbf{v}}{\partial x} + \left(\frac{\partial \mathbf{u}}{\partial x} \right)^T \mathbf{v}$ <p>This is the derivative of scalar product of two vectors when $\frac{\partial \mathbf{u}}{\partial x}$ and $\frac{\partial \mathbf{v}}{\partial x}$ are tangent vectors.</p>
$\mathbf{U} = \mathbf{U}(x)$ is matrix-by-scalar function	$\frac{\partial \mathbf{U} }{\partial x} = \mathbf{U} \text{tr} \left(\mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \right)$
$\mathbf{U} = \mathbf{U}(x)$ is matrix-by-scalar function	$\frac{\partial \ln \mathbf{U} }{\partial x} = \text{tr} \left(\mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \right)$
$\mathbf{U} = \mathbf{U}(x)$ is matrix-by-scalar function	$\frac{\partial^2 \mathbf{U} }{\partial x^2} = \mathbf{U} \left(\text{tr} \left(\mathbf{U}^{-1} \frac{\partial^2 \mathbf{U}}{\partial x^2} \right) + \left(\text{tr} \left(\mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \right) \right)^2 - \text{tr} \left(\left(\mathbf{U}^{-1} \frac{\partial \mathbf{U}}{\partial x} \right)^2 \right) \right)$
$\mathbf{U} = \mathbf{U}(x)$ is matrix-by-scalar function. f is scalar-by-matrix function	$\frac{\partial g(\mathbf{U})}{\partial x} = \text{tr} \left(\frac{\partial g(\mathbf{U})}{\partial \mathbf{U}} \frac{\partial \mathbf{U}}{\partial x} \right)$
A is square matrix constant. $\mathbf{G}(\cdot)$ is any matrix polynomial $\mathbf{P}(X)$ or matrix-power-series function such as $\mathbf{E}^X, \mathbf{LNX}, \mathbf{SINX}, \mathbf{COSX}$.	$\frac{\partial \text{tr}(\mathbf{G}(xA))}{\partial x} = \text{tr}(A \mathbf{G}'(xA))$ <p>Following are some popular derivatives</p> $\begin{aligned} \frac{\partial \text{tr}(\mathbf{P}(xA))}{\partial x} &= \frac{\partial \text{tr}(\sum_{k=0}^n a_k(xA)^k)}{\partial x} = \text{tr}(A \mathbf{P}'(xA)) \\ &= \text{tr} \left(A \sum_{k=1}^n k a_k(xA)^{k-1} \right) \end{aligned}$

$$\frac{\partial \text{tr}((xA)^n)}{\partial x} = \text{tr}(nA(xA)^{n-1}) \text{ where } n \text{ is positive integer}$$

$$\frac{\partial \text{tr}(E^{xA})}{\partial x} = \text{tr}(AE^{xA})$$

$$\frac{\partial \text{tr}(\mathbf{SIN}(xA))}{\partial x} = \text{tr}(A\mathbf{COS}(xA))$$

$$\frac{\partial \text{tr}(\mathbf{COS}(xA))}{\partial x} = -\text{tr}(A\mathbf{SIN}(xA))$$

5. Some applications of matrix analysis and calculus

There are a lot of applications of matrix analysis and calculus, for example, statistics, data analysis, optimization, geometry, kinematics, and dynamics. These applications are huge and there are many subjects related them. In the short report, we only introduce typical examples of matrix analysis and calculus. Some aforementioned examples such as Jordan decomposition and singular value decomposition are also good examples for interesting matrix applications. This section has two examples:

- Euclidean distance with general metric matrix (Härdle & Simar, 2013, pp. 71-74) illustrates applications of matrix analysis.
- Optimization with Lagrangian duality illustrates applications of matrix calculus.

Euclidean distance with general metric matrix

Distance d is defined as the function from $\mathbb{R}^n \times \mathbb{R}^n$ to \mathbb{R}_+ :

$$\mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}_+$$

$$d: (\mathbf{x}, \mathbf{y}) \rightarrow r \geq 0$$

Where \mathbf{x} and \mathbf{y} are vectors in vector space over real field \mathbb{R}^n .

Distance d is always greater than or equal to 0, which satisfies three following axioms (Härdle & Simar, 2013, p. 71).

$$\begin{cases} d(\mathbf{x}, \mathbf{y}) > 0 \quad \forall \mathbf{x} \neq \mathbf{y} \\ d(\mathbf{x}, \mathbf{y}) = 0 \text{ if and only if } \mathbf{x} = \mathbf{y} \\ d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}) \\ d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z}) + d(\mathbf{y}, \mathbf{z}) \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \end{cases}$$

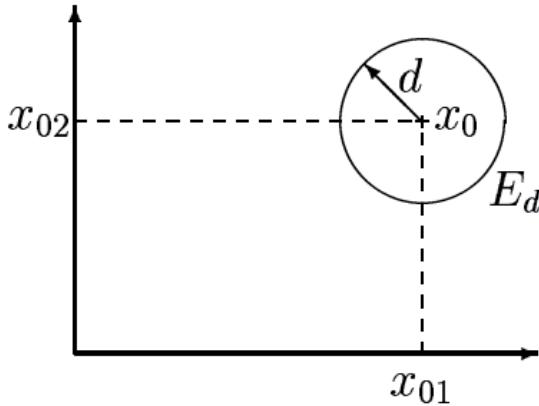
The Euclidean distance between two points is defined as following equation (Härdle & Simar, 2013, p. 71).

$$d^2(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{A} (\mathbf{x} - \mathbf{y})$$

Where \mathbf{A} is positive definite matrix ($\mathbf{A} > 0$) and \mathbf{A} is called a metric and the space on which \mathbf{A} is defined is called metric space. Euclidean distance is concerned. If \mathbf{A} is identical matrix $\mathbf{A} = \mathbf{I}_n$, Euclidean distance defined by $d^2(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n (x_i - y_i)^2$. Given a point \mathbf{x}_0 and a scalar constant d , a n -dimension sphere with radius d is defined as following equation.

$$(\mathbf{x} - \mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) = d^2$$

Following is figure of d -radius sphere (Härdle & Simar, 2013, p. 72):

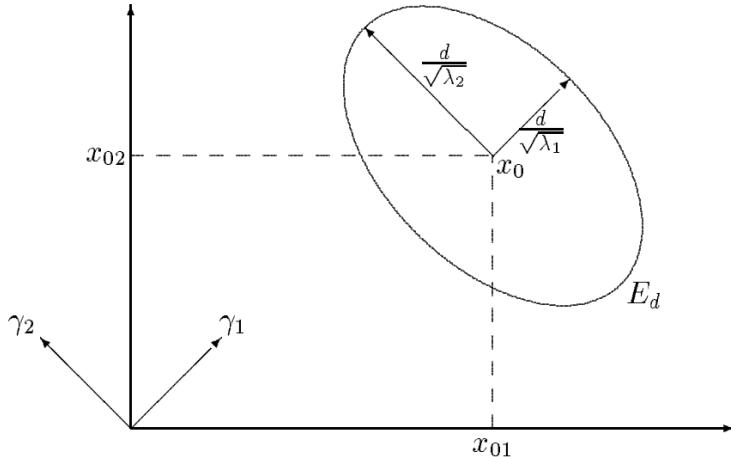


Where E_d denotes Euclidean metric space.

This sphere is a set of points which are far from x_0 a distance d . Given a point x_0 , a matrix $A(n \times n)$ and a scalar d , a n -dimension ellipsoid is defined as following equation.

$$(x - x_0)^T A (x - x_0) = d^2$$

Following is figure of Ellipsoid with center x_0 , matrix A and constant d (Härdle & Simar, 2013, p. 72):



Note that A is invertible matrix. Suppose matrix A has n eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ and n respective orthogonal eigenvectors $\gamma_1, \gamma_2, \dots, \gamma_n$. This ellipsoid has following properties (Härdle & Simar, 2013, p. 73):

- The principle axes of ellipsoid have the same direction to eigenvectors $\gamma_1, \gamma_2, \dots, \gamma_n$. Of course, the number of axes is equal to the number of eigenvectors.

- The half of length of each axes is equal to $\sqrt{\frac{d^2}{\lambda_i}}$ where λ_i (s) are eigenvalues.
- Let x_{0i} be the element i of center $\mathbf{x}_0 = \{x_{01}, x_{02}, \dots, x_{0n}\}$. The n -dimension rectangle surrounding ellipsoid is determined by equation: $x_{0i} - \sqrt{d^2 a^{ii}} \leq x_i \leq x_{0i} + \sqrt{d^2 a^{ii}}$ where a^{ii} is the element (i, i) of A^{-1} .
- The coordinate of tangency point (x_i) between ellipsoid and its surrounding rectangle in the positive direction of j^{th} axis is $x_i = \left(a^{ij} \sqrt{\frac{d^2}{a^{jj}}} \right)$ where a^{ij} and a^{jj} is the elements (i, j) and (j, j) of A^{-1} .

Optimization with Lagrangian function

Given a scalar-by-vector function $f(\mathbf{x})$, there is a requirement “how to minimize f with a scalar-by-vector constraint $h(\mathbf{x})$ ” This is optimization problem also called non-linear programming, specialized by following extreme-value expression.

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & h(\mathbf{x}) = 0 \end{array}$$

Where $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n$ and f, h have continuous second partial derivatives. Function f is called target function and h is called constraint function.

Suppose \mathbf{x}^* is extreme point of target function f satisfying constraint $h(\mathbf{x}^*) = 0$. Given the hyper-surface (surface in \mathbb{R}^3) determined equation $h(\mathbf{x}) = 0$, for each curve $\gamma(t)$ lying on this hyper-surface together with condition that $\gamma(t)$ goes through \mathbf{x}^* at $t = 0$, we have (Dinh, Pham, & Ta, 2002, pp. 59-61):

$$\gamma: \mathbb{R} \rightarrow \mathbb{R}^n$$

$$\gamma(0) = \mathbf{x}^*$$

$$h(\gamma(0)) = h(\mathbf{x}^*) = 0$$

Note that $\gamma(t)$ is vector-by-scalar function. The tangent vector of $\gamma(t)$ denoted τ at point \mathbf{x}^* is:

$$\tau = \gamma'(0) = \frac{\partial \gamma}{\partial t}(0)$$

The tangent space (tangent plane in \mathbb{R}^3) denoted \mathbb{T} of hyper-surface $h(\mathbf{x}) = 0$ at point \mathbf{x}^* is composed of all tangent vectors of curves $\gamma(t)$.

$$\mathbb{T} = \{\tau = \gamma'(0) \text{ for all curves } \gamma(t)\}$$

The derivative of $h(\mathbf{x})$ which is the gradient ∇h is orthogonal to tangent space \mathbb{T} at \mathbf{x}^* ; in other words it lies in the orthogonal complement \mathbb{T}^\perp of tangent space \mathbb{T} .

$$\nabla h(\mathbf{x}^*)\boldsymbol{\tau} = 0, \forall \boldsymbol{\tau} \in \mathbb{T}$$

Function f gets extreme value at \mathbf{x}^* and so its first-order derivative at \mathbf{x}^* is 0. According to chain rule (Dinh, Pham, & Ta, 2002, pp. 49-50), we have:

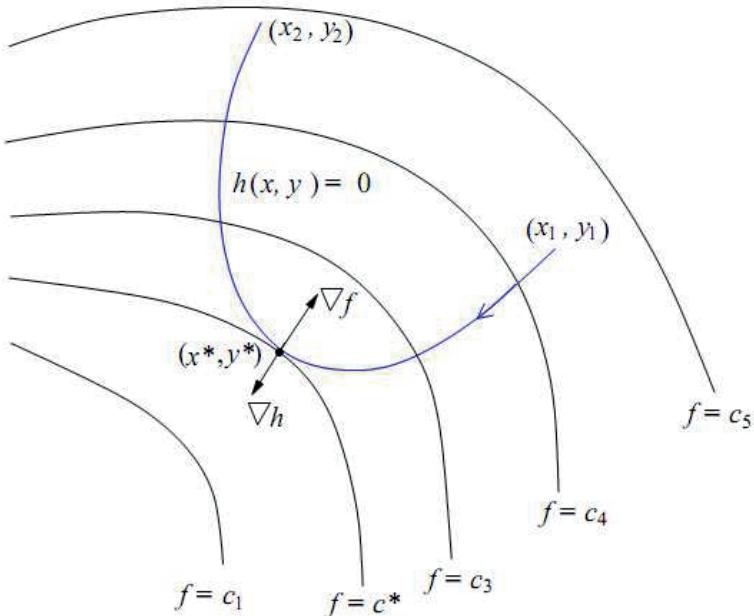
$$0 = f'(\mathbf{x}^*) = f'(\gamma(0)) = \nabla f(\mathbf{x}^*)\gamma'(0) = \nabla f(\mathbf{x}^*)\boldsymbol{\tau}, \forall \boldsymbol{\tau} \in \mathbb{T}$$

Where ∇f is gradient of f .

Hence, both gradients ∇f and ∇h are orthogonal to all tangent vector $\boldsymbol{\tau}$ at extreme point \mathbf{x}^* , which leads to conclusion that gradient ∇f is parallel with the gradient ∇h at \mathbf{x}^* . Consequently, it is reasoned out that there exists a scalar λ such that:

$$\nabla f(\mathbf{x}^*) + \lambda \nabla h(\mathbf{x}^*) = 0$$

For illustrating the assertion “gradient ∇f is parallel with gradient ∇h at extreme point \mathbf{x}^* ”, we suppose that $f(x, y)$ and $h(x, y)$ are 2-variable functions. The constraint function $h(x, y)$ is reduced to a curve superposed by a set of contours of target function f . Suppose we have 5 contours $f = c_1 < f = c^* < f = c_3 < f = c_4 < f = c_5$ where each value c_i is the value of contour of f and c^* is extreme value of f . Following is the figure (Jia, 2013, p. 2) illustrating constraint function $h(x, y)$ and contours of target function $f(x, y)$.



Please pay attention that when projecting the gradient vector of function $y = f(x, y)$ onto the plane containing its contour $f(x, y) = c$ then the gradient is orthogonal to the tangent vector of the contour. According to (Jia, 2013, pp. 2-3), it is easy to recognize that when moving along the curve $h(x, y) = 0$, the value of contour is increased and decreased and function f gets extreme value so that $h(x, y) =$

0 at point (x^*, y^*) . Although $f = c_1$ is minimum but $h(x, y)$ does not intersect with contour $f(x, y) = c_1$. It means that gradient ∇f is parallel with gradient ∇h .

In general, the extreme points x^* are solutions of equations:

$$\begin{cases} h(x) = 0 \\ \nabla f + \lambda \nabla h = 0 \end{cases}$$

Let $l(x, \lambda) = f(x) + \lambda h(x)$ be Lagrangian scalar-by-vector function, we have:

$$\frac{\partial l(x, \lambda)}{\partial (x, \lambda)} = \nabla l = (\nabla f + \lambda \nabla h, h)$$

Note that both x and λ are independent variable in Lagrangian function and so the notation $\frac{\partial l(x, \lambda)}{\partial (x, \lambda)}$ indicates derivative of Lagrangian function w.r.t vector $(x^T, \lambda) = (x_1, x_2, \dots, x_n, \lambda)$.

Derived from conclusion “extreme points x^* are solutions of equations $h(x) = 0$ and $\nabla f + \lambda \nabla h = 0$ ”, principle Lagrange states that if x^* is extreme point, it always exists the real scalar λ^* so that the gradient ∇l is equal to $\mathbf{0}$ when substituting (x^*, λ^*) into it.

$$\frac{\partial l(x, \lambda)}{\partial (x, \lambda)}(x^*, \lambda^*) = \nabla L(x^*, \lambda^*) = \begin{pmatrix} \nabla f(x^*) + \lambda^* \nabla h(x^*) \\ h(x^*) \end{pmatrix}^T = (\mathbf{0}_n^T, 0)$$

Note that $(\mathbf{0}_n^T, 0)$ is zero $1_{x(n+1)}$ row vector. In brief, we have:

$$\nabla l(x^*, \lambda^*) = \mathbf{0}_{n+1}^T$$

Note that Lagrangian function is scalar-by-vector function and so its derivative is really row gradient vector. Variables x^* and λ^* are always co-existent, hence, λ is called Lagrangian multiplier or Lagrangian dual variable. Lagrange principle is known as Lagrangian duality. If $x \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}$ then, Lagrangian function has $n + 1$ partial elements and so it has $n + 1$ partial derivatives. Therefore, in practice, we set $n + 1$ these partial derivatives to be 0 so as to compose a set of $n + 1$ equations. After that Lagrangian multiplier λ^* and extreme point x^* are solution of such equations, as follows:

$$\begin{cases} \frac{\partial f}{\partial x_1} + \lambda \frac{\partial h}{\partial x_1} = 0 \\ \frac{\partial f}{\partial x_2} + \lambda \frac{\partial h}{\partial x_2} = 0 \\ \vdots \\ \frac{\partial f}{\partial x_n} + \lambda \frac{\partial h}{\partial x_n} = 0 \\ h(x) = 0 \end{cases}$$

If it is too difficult to solve these equations by primary transformation, there are many methods to solve them, for example, Newton's method.

For example, let us find out extreme value of $f(x, y) = xy$ with constraint

$$h(x, y) = x^2/8 + y^2/2 - 1 = 0$$

The Lagrangian function and its gradient are constructed as below:

$$l(x, y, \lambda) = xy + \lambda(x^2/8 + y^2/2 - 1)$$

$$\nabla l = \frac{\partial l(\mathbf{x}, \lambda)}{\partial (\mathbf{x}, \lambda)} = \left(y + \frac{\lambda x}{4}, x + \lambda y, \frac{x^2}{8} + \frac{y^2}{2} - 1 \right)$$

Lagrangian multiplier λ^* and extreme point \mathbf{x}^* are solutions of equations constructed from setting ∇l to be $\mathbf{0}$.

$$\nabla l = \mathbf{0} \Leftrightarrow \begin{cases} y + \frac{\lambda x}{4} = 0 \\ x + \lambda y = 0 \\ \frac{x^2}{8} + \frac{y^2}{2} - 1 = 0 \end{cases}$$

$$\Rightarrow \begin{cases} \lambda = 2, \text{ two extreme points } (x = -2, y = 1) \text{ or } (x = 2, y = -1) \\ \lambda = -2, \text{ two extreme points } (x = 2, y = 1) \text{ or } (x = -2, y = -1) \end{cases}$$

We have:

$$\begin{cases} f(-2, 1) = -2 \\ f(2, -1) = -2 \\ f(2, 1) = 2 \\ f(-2, -1) = 2 \end{cases}$$

Therefore, f gets local maximum at $(2, 1), (-2, -1)$ and local minimum at $(-2, 1), (2, -1)$.

Now the optimization problem is extended with multi-constraints

$$\begin{array}{ll} \text{Minimize} & f(\mathbf{x}) \\ \text{subject to} & h_1(\mathbf{x}) = 0 \\ & h_2(\mathbf{x}) = 0 \\ & \vdots \\ & h_m(\mathbf{x}) = 0 \end{array}$$

Or

$$\begin{array}{ll} \text{Minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{h}(\mathbf{x}) = \mathbf{0} \end{array}$$

Hence, the constraint function becomes vector-by-vector function $\mathbf{h}(\mathbf{x}) = \begin{pmatrix} h_1(\mathbf{x}) \\ h_2(\mathbf{x}) \\ \vdots \\ h_m(\mathbf{x}) \end{pmatrix}$. For convenience, the derivative of function $\mathbf{h}(\mathbf{x})$ which is Jacobian matrix is denoted by gradient notation.

$$\nabla \mathbf{h} = \frac{\partial \mathbf{h}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial h_1}{\partial \mathbf{x}} \\ \frac{\partial h_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial h_m}{\partial \mathbf{x}} \end{pmatrix} = \begin{pmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \cdots & \frac{\partial h_1}{\partial x_n} \\ \frac{\partial h_2}{\partial x_1} & \frac{\partial h_2}{\partial x_2} & \cdots & \frac{\partial h_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_m}{\partial x_1} & \frac{\partial h_m}{\partial x_2} & \cdots & \frac{\partial h_m}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \nabla h_1 \\ \nabla h_2 \\ \vdots \\ \nabla h_m \end{pmatrix}$$

$$\text{Where } \nabla h_i = \left(\frac{\partial h_i}{\partial x_1}, \frac{\partial h_i}{\partial x_2}, \dots, \frac{\partial h_i}{\partial x_n} \right), \forall i = \overline{1, m}$$

Given extreme point \mathbf{x}^* satisfying $\mathbf{h}(\mathbf{x}^*) = \mathbf{0}$ is regular point if Jacobian matrix $\nabla \mathbf{h}$ is non-singular at \mathbf{x}^* . The tangent space \mathbb{T} at \mathbf{x}^* exists if and only if \mathbf{x}^* is regular point and so the Lagrangian function exists if and only if Jacobian matrix $\nabla \mathbf{h}$ is non-singular at \mathbf{x}^* . The tangent space at \mathbf{x}^* is null space of $\nabla \mathbf{h}(\mathbf{x}^*)$, it means:

$$\forall \boldsymbol{\tau} \in \mathbb{T}, \nabla \mathbf{h}(\mathbf{x}^*) \boldsymbol{\tau} = \mathbf{0}$$

Note that $\mathbf{0}$ is zero column vector and we still have:

$$\nabla f(\mathbf{x}^*) \boldsymbol{\tau} = 0, \forall \boldsymbol{\tau} \in \mathbb{T}$$

Both gradient ∇f and Jacobian matrix $\nabla \mathbf{h}$ lie in complement \mathbb{T}^\perp of tangent space \mathbb{T} at \mathbf{x}^* , so there is

a real number vector $\boldsymbol{\lambda} = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_m \end{pmatrix}$ such that:

$$\nabla f(\mathbf{x}^*) + \boldsymbol{\lambda}^T \nabla \mathbf{h}(\mathbf{x}^*) = \mathbf{0}$$

and

$$\mathbf{h}(\mathbf{x}^*) = \mathbf{0}$$

Lagrangian function is re-constructed:

$$l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x})$$

According to principle Lagrange, if \mathbf{x}^* is extreme point, it always exists the $m \times 1$ vector $\boldsymbol{\lambda}^*$ such that

$$\nabla l(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \begin{pmatrix} \nabla f(\mathbf{x}^*) + \boldsymbol{\lambda}^{*T} \nabla \mathbf{h}(\mathbf{x}^*) \\ (\mathbf{h}(\mathbf{x}^*))^T \end{pmatrix}^T = (\mathbf{0}_n^T, \mathbf{0}_m^T)$$

Note that $(\mathbf{0}_n^T, \mathbf{0}_m^T)$ is zero $1 \times (n+m)$ row vector. In brief, we have:

$$\nabla l(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}_{n+m}^T$$

Lagrangian function with m constraints has $n + m$ partial elements and so it has $n + m$ partial derivatives. Therefore, in practice, we set $n + m$ these partial derivatives to be 0 so as to construct a set of $n + m$ equations whose solutions are extreme point \mathbf{x}^* and Lagrangian multipliers vector $\boldsymbol{\lambda}^*$, as follows:

$$\begin{cases} \frac{\partial f}{\partial x_1} + \sum_{i=1}^m \lambda_i \frac{\partial h_i}{\partial x_1} = 0 \\ \frac{\partial f}{\partial x_2} + \sum_{i=1}^m \lambda_i \frac{\partial h_i}{\partial x_2} = 0 \\ \vdots \\ \frac{\partial f}{\partial x_n} + \sum_{i=1}^m \lambda_i \frac{\partial h_i}{\partial x_n} = 0 \\ h_1(\mathbf{x}) = 0 \\ h_2(\mathbf{x}) = 0 \\ \vdots \\ h_m(\mathbf{x}) = 0 \end{cases}$$

Lagrangian duality can be extended to solve optimization problem with inequality constraints. This is the most general case of non-linear programming, as follows:

$$\text{Minimize } f(\mathbf{x})$$

$$\text{subject to } \begin{aligned} \mathbf{h}(\mathbf{x}) &= \mathbf{0} \\ \mathbf{g}(\mathbf{x}) &\leq \mathbf{0} \end{aligned}$$

Where $\mathbf{h}(\mathbf{x}) = \begin{pmatrix} h_1(\mathbf{x}) \\ h_2(\mathbf{x}) \\ \vdots \\ h_m(\mathbf{x}) \end{pmatrix}$ is equality constraint vector-by-vector function and $\mathbf{g}(\mathbf{x}) = \begin{pmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \\ \vdots \\ g_p(\mathbf{x}) \end{pmatrix}$ is

inequality constraint vector-by-vector function. Lagrangian function is re-constructed:

$$l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x})$$

Where $\boldsymbol{\lambda} = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_m \end{pmatrix}$ is real number vector and $\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix}$ with attention that $\mu_i \geq 0, \forall i$. According

to principle Lagrange, if \mathbf{x}^* is extreme point, it always exists $m \times 1$ vector $\boldsymbol{\lambda}^*$ and $p \times 1$ vector $\boldsymbol{\mu}^*$ such that

$$\nabla l(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \begin{pmatrix} \nabla f(\mathbf{x}^*) + \boldsymbol{\lambda}^{*T} \nabla \mathbf{h}(\mathbf{x}^*) + \boldsymbol{\mu}^{*T} \nabla \mathbf{g}(\mathbf{x}^*) \\ (\mathbf{h}(\mathbf{x}^*))^T \\ (\mathbf{g}(\mathbf{x}^*))^T \end{pmatrix}^T = (\mathbf{0}_n^T, \mathbf{0}_m^T, \mathbf{0}_p^T)$$

Where,

$$\nabla \mathbf{h} = \frac{\partial \mathbf{h}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial h_1}{\partial \mathbf{x}} \\ \frac{\partial h_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial h_m}{\partial \mathbf{x}} \end{pmatrix} = \begin{pmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \cdots & \frac{\partial h_1}{\partial x_n} \\ \frac{\partial h_2}{\partial x_1} & \frac{\partial h_2}{\partial x_2} & \cdots & \frac{\partial h_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_m}{\partial x_1} & \frac{\partial h_m}{\partial x_2} & \cdots & \frac{\partial h_m}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \nabla h_1 \\ \nabla h_2 \\ \vdots \\ \nabla h_m \end{pmatrix}$$

$$\nabla h_i = \left(\frac{\partial h_i}{\partial x_1}, \frac{\partial h_i}{\partial x_2}, \dots, \frac{\partial h_i}{\partial x_n} \right), \forall i = \overline{1, m}$$

$$\nabla \mathbf{g} = \frac{\partial \mathbf{g}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial g_1}{\partial \mathbf{x}} \\ \frac{\partial g_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial g_p}{\partial \mathbf{x}} \end{pmatrix} = \begin{pmatrix} \frac{\partial g_1}{\partial x_1} & \frac{\partial g_1}{\partial x_2} & \cdots & \frac{\partial g_1}{\partial x_n} \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} & \cdots & \frac{\partial g_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g_p}{\partial x_1} & \frac{\partial g_p}{\partial x_2} & \cdots & \frac{\partial g_p}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \nabla g_1 \\ \nabla g_2 \\ \vdots \\ \nabla g_p \end{pmatrix}$$

$$\nabla g_i = \left(\frac{\partial g_i}{\partial x_1}, \frac{\partial g_i}{\partial x_2}, \dots, \frac{\partial g_i}{\partial x_n} \right), \forall i = \overline{1, p}$$

Note that $(\mathbf{0}_n^T, \mathbf{0}_m^T, \mathbf{0}_p^T)$ is zero $1_{x(n+m+p)}$ row vector. In brief, we have:

$$\nabla l(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \mathbf{0}_{n+m+p}^T$$

In practice, we set $n + m + p$ these partial derivatives to be 0 so as to construct a set of $n + m$ equations whose solutions are extreme point \mathbf{x}^* and Lagrangian multipliers vectors $\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*$ as follows:

$$\left\{ \begin{array}{l} \frac{\partial f}{\partial x_1} + \sum_{i=1}^m \lambda_i \frac{\partial h_i}{\partial x_1} + \sum_{j=1}^p \mu_j \frac{\partial g_j}{\partial x_1} = 0 \\ \frac{\partial f}{\partial x_2} + \sum_{i=1}^m \lambda_i \frac{\partial h_i}{\partial x_2} + \sum_{j=1}^p \mu_j \frac{\partial g_j}{\partial x_2} = 0 \\ \vdots \\ \frac{\partial f}{\partial x_n} + \sum_{i=1}^m \lambda_i \frac{\partial h_i}{\partial x_n} + \sum_{j=1}^p \mu_j \frac{\partial g_j}{\partial x_n} = 0 \\ h_1(\mathbf{x}) = 0 \\ h_2(\mathbf{x}) = 0 \\ \vdots \\ h_m(\mathbf{x}) = 0 \\ g_1(\mathbf{x}) = 0 \\ g_2(\mathbf{x}) = 0 \\ \vdots \\ g_p(\mathbf{x}) = 0 \end{array} \right.$$

In general, we need to solve these equations above to find out optimal solutions of optimization problem.

6. Conclusion

Now we researched main subjects of matrix analysis and calculus over the report in which sections “Matrix analysis”, “Matrix derivative” and “Composite derivative” are important ones. Matrix analysis focuses on data analyzing and data processing whereas matrix calculus focuses on real-time process, differential geometry, and optimization. Some domains apply both matrix analysis and matrix calculus. For example, multivariate statistics uses not only data analysis techniques but also other studies such as differential, calculus, graphics, combinatorics, and probability for analyzing data and testing hypothesis. Matrix analysis is very necessary to matrix calculus and so we should research it before studying matrix calculus and differential. Some enhanced subjects such as matrix-by-matrix functions and derivatives, Kronecker product, Hadamard product, complex matrix, and multivariate distribution are not mentioned in the report. We will discuss them in next research.

Acknowledgement

The report expresses my deep gratitude to Professor Ho, Minh T. at Vietnam Institute of Mathematics who taught me a lot of valuable knowledge and gave me the motivation to complete both our knowledge and the report.

Bibliography

- Baker, K. (2013). *Singular Value Decomposition Tutorial*. The College of Humanities and Social Sciences, Linguistics Department Advisory Board. Montclair State University.
- Baker, M. J. (n.d.). *Maths - Powers of Vectors*. Retrieved 2014, from EuclideanSpace - Mathematics and Computing: <http://www.euclideanspace.com/maths/algebra/vectors/vecAlgebra/powers/index.htm>
- Dinh, L. T., Pham, D. H., & Ta, P. D. (2002). *Multivariate Analysis - Principles and Practices* (Vol. II). (K. H. Ha, T. V. Ngo, & D. H. Pham, Eds.) Hanoi, Vietnam: Hanoi National University Publisher.
- Härdle, W., & Simar, L. (2013). *Applied Multivariate Statistical Analysis*. Berlin, Germany: Research Data Center, School of Business and Economics, Humboldt University.
- Hoang, V. H. (2012, December 24). *Linear mapping - Diagonalizing matrix*. (V. H. Hoang, Performer) YouTube, Ho Chi Minh, Vietnam. Retrieved from <https://www.youtube.com/watch?v=NkSDymM-qPg>
- Jia, Y.-B. (2013). *Lagrange Multipliers*. Lecture notes on course “Problem Solving Techniques for Applied Computer Science”, Iowa State University of Science and Technology, USA.
- Nguyen, V. H. (1999). *Linear Algebra*. Hanoi, Vietnam: Hanoi National University Publishing House.
- Petersen, K. B., & Pedersen, M. S. (2012). *The Matrix Cookbook*. Copenhagen, Denmark: SAS Institute. Retrieved from <http://orion.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>
- Wikipedia. (2014, March 9). *Cayley–Hamilton theorem*. (Wikimedia Foundation) Retrieved from Wikipedia website: http://en.wikipedia.org/wiki/Cayley%E2%80%93Hamilton_theorem
- Wikipedia. (2014, February 25). *Logarithm of a matrix*. (Wikimedia Foundation) Retrieved from Wikipedia website: http://en.wikipedia.org/wiki/Matrix_logarithm
- Wikipedia. (2014, March 3). *Matrix calculus*. (Wikimedia Foundation) Retrieved from Wikipedia website: http://en.wikipedia.org/wiki/Matrix_calculus
- Wikipedia. (2014, February 20). *Matrix exponential*. (Wikimedia Foundation) Retrieved from Wikipedia website: http://en.wikipedia.org/wiki/Matrix_exponential

Wikipedia. (2014, February 12). *Matrix function*. (Wikimedia Foundation) Retrieved from Wikipedia website: http://en.wikipedia.org/wiki/Matrix_function



yes I want morebooks!

Buy your books fast and straightforward online - at one of the world's fastest growing online book stores! Environmentally sound due to Print-on-Demand technologies.

Buy your books online at
www.get-morebooks.com

Kaufen Sie Ihre Bücher schnell und unkompliziert online – auf einer der am schnellsten wachsenden Buchhandelsplattformen weltweit!
Dank Print-On-Demand umwelt- und ressourcenschonend produziert.

Bücher schneller online kaufen
www.morebooks.de

OmniScriptum Marketing DEU GmbH
Heinrich-Böcking-Str. 6-8
D - 66121 Saarbrücken
Telefax: +49 681 93 81 567-9

info@omniscriptum.com
www.omniscriptum.com



