

The Matrix Cookbook

[<http://matrixcookbook.com>]

Kaare Brandt Petersen
Michael Syskind Pedersen

VERSION: NOVEMBER 15, 2012

Introduction

What is this? These pages are a collection of facts (identities, approximations, inequalities, relations, ...) about matrices and matters relating to them. It is collected in this form for the convenience of anyone who wants a quick desktop reference .

Disclaimer: The identities, approximations and relations presented here were obviously not invented but collected, borrowed and copied from a large amount of sources. These sources include similar but shorter notes found on the internet and appendices in books - see the references for a full list.

Errors: Very likely there are errors, typos, and mistakes for which we apologize and would be grateful to receive corrections at cookbook@2302.dk.

Its ongoing: The project of keeping a large repository of relations involving matrices is naturally ongoing and the version will be apparent from the date in the header.

Suggestions: Your suggestion for additional content or elaboration of some topics is most welcome acookbook@2302.dk.

Keywords: Matrix algebra, matrix relations, matrix identities, derivative of determinant, derivative of inverse matrix, differentiate a matrix.

Acknowledgements: We would like to thank the following for contributions and suggestions: Bill Baxter, Brian Templeton, Christian Rishøj, Christian Schröppel, Dan Boley, Douglas L. Theobald, Esben Hoegh-Rasmussen, Evripidis Karseras, Georg Martius, Glynne Casteel, Jan Larsen, Jun Bin Gao, Jürgen Struckmeier, Kamil Dedecius, Karim T. Abou-Moustafa, Korbinian Strimmer, Lars Christiansen, Lars Kai Hansen, Leland Wilkinson, Ligu He, Loic Thibaut, Markus Froeb, Michael Hubatka, Miguel Barão, Ole Winther, Pavel Sakov, Stephan Hattinger, Troels Pedersen, Vasile Sima, Vincent Rabaud, Zhaoshui He. We would also like thank The Oticon Foundation for funding our PhD studies.

Contents

1	Basics	6
1.1	Trace	6
1.2	Determinant	6
1.3	The Special Case 2x2	7
2	Derivatives	8
2.1	Derivatives of a Determinant	8
2.2	Derivatives of an Inverse	9
2.3	Derivatives of Eigenvalues	10
2.4	Derivatives of Matrices, Vectors and Scalar Forms	10
2.5	Derivatives of Traces	12
2.6	Derivatives of vector norms	14
2.7	Derivatives of matrix norms	14
2.8	Derivatives of Structured Matrices	14
3	Inverses	17
3.1	Basic	17
3.2	Exact Relations	18
3.3	Implication on Inverses	20
3.4	Approximations	20
3.5	Generalized Inverse	21
3.6	Pseudo Inverse	21
4	Complex Matrices	24
4.1	Complex Derivatives	24
4.2	Higher order and non-linear derivatives	26
4.3	Inverse of complex sum	27
5	Solutions and Decompositions	28
5.1	Solutions to linear equations	28
5.2	Eigenvalues and Eigenvectors	30
5.3	Singular Value Decomposition	31
5.4	Triangular Decomposition	32
5.5	LU decomposition	32
5.6	LDM decomposition	33
5.7	LDL decompositions	33
6	Statistics and Probability	34
6.1	Definition of Moments	34
6.2	Expectation of Linear Combinations	35
6.3	Weighted Scalar Variable	36
7	Multivariate Distributions	37
7.1	Cauchy	37
7.2	Dirichlet	37
7.3	Normal	37
7.4	Normal-Inverse Gamma	37
7.5	Gaussian	37
7.6	Multinomial	37

7.7	Student's t	37
7.8	Wishart	38
7.9	Wishart, Inverse	39
8	Gaussians	40
8.1	Basics	40
8.2	Moments	42
8.3	Miscellaneous	44
8.4	Mixture of Gaussians	44
9	Special Matrices	46
9.1	Block matrices	46
9.2	Discrete Fourier Transform Matrix, The	47
9.3	Hermitian Matrices and skew-Hermitian	48
9.4	Idempotent Matrices	49
9.5	Orthogonal matrices	49
9.6	Positive Definite and Semi-definite Matrices	50
9.7	Singleentry Matrix, The	52
9.8	Symmetric, Skew-symmetric/Antisymmetric	54
9.9	Toeplitz Matrices	54
9.10	Transition matrices	55
9.11	Units, Permutation and Shift	56
9.12	Vandermonde Matrices	57
10	Functions and Operators	58
10.1	Functions and Series	58
10.2	Kronecker and Vec Operator	59
10.3	Vector Norms	61
10.4	Matrix Norms	61
10.5	Rank	62
10.6	Integral Involving Dirac Delta Functions	62
10.7	Miscellaneous	63
A	One-dimensional Results	64
A.1	Gaussian	64
A.2	One Dimensional Mixture of Gaussians	65
B	Proofs and Details	66
B.1	Misc Proofs	66

Notation and Nomenclature

\mathbf{A}	Matrix
\mathbf{A}_{ij}	Matrix indexed for some purpose
\mathbf{A}_i	Matrix indexed for some purpose
\mathbf{A}^{ij}	Matrix indexed for some purpose
\mathbf{A}^n	Matrix indexed for some purpose or The n .th power of a square matrix
\mathbf{A}^{-1}	The inverse matrix of the matrix \mathbf{A}
\mathbf{A}^+	The pseudo inverse matrix of the matrix \mathbf{A} (see Sec. 3.6)
$\mathbf{A}^{1/2}$	The square root of a matrix (if unique), not elementwise
$(\mathbf{A})_{ij}$	The (i, j) .th entry of the matrix \mathbf{A}
A_{ij}	The (i, j) .th entry of the matrix \mathbf{A}
$[\mathbf{A}]_{ij}$	The ij -submatrix, i.e. \mathbf{A} with i .th row and j .th column deleted
\mathbf{a}	Vector (column-vector)
\mathbf{a}_i	Vector indexed for some purpose
a_i	The i .th element of the vector \mathbf{a}
a	Scalar
$\Re z$	Real part of a scalar
$\Re \mathbf{z}$	Real part of a vector
$\Re \mathbf{Z}$	Real part of a matrix
$\Im z$	Imaginary part of a scalar
$\Im \mathbf{z}$	Imaginary part of a vector
$\Im \mathbf{Z}$	Imaginary part of a matrix
$\det(\mathbf{A})$	Determinant of \mathbf{A}
$\text{Tr}(\mathbf{A})$	Trace of the matrix \mathbf{A}
$\text{diag}(\mathbf{A})$	Diagonal matrix of the matrix \mathbf{A} , i.e. $(\text{diag}(\mathbf{A}))_{ij} = \delta_{ij} A_{ij}$
$\text{eig}(\mathbf{A})$	Eigenvalues of the matrix \mathbf{A}
$\text{vec}(\mathbf{A})$	The vector-version of the matrix \mathbf{A} (see Sec. 10.2.2)
\sup	Supremum of a set
$\ \mathbf{A}\ $	Matrix norm (subscript if any denotes what norm)
\mathbf{A}^T	Transposed matrix
\mathbf{A}^{-T}	The inverse of the transposed and vice versa, $\mathbf{A}^{-T} = (\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$.
\mathbf{A}^*	Complex conjugated matrix
\mathbf{A}^H	Transposed and complex conjugated matrix (Hermitian)
$\mathbf{A} \circ \mathbf{B}$	Hadamard (elementwise) product
$\mathbf{A} \otimes \mathbf{B}$	Kronecker product
$\mathbf{0}$	The null matrix. Zero in all entries.
\mathbf{I}	The identity matrix
\mathbf{J}^{ij}	The single-entry matrix, 1 at (i, j) and zero elsewhere
Σ	A positive definite matrix
Λ	A diagonal matrix

1 Basics

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

$$(\mathbf{ABC}\dots)^{-1} = \dots\mathbf{C}^{-1}\mathbf{B}^{-1}\mathbf{A}^{-1} \quad (2)$$

$$(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T \quad (3)$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T \quad (4)$$

$$(\mathbf{AB})^T = \mathbf{B}^T\mathbf{A}^T \quad (5)$$

$$(\mathbf{ABC}\dots)^T = \dots\mathbf{C}^T\mathbf{B}^T\mathbf{A}^T \quad (6)$$

$$(\mathbf{A}^H)^{-1} = (\mathbf{A}^{-1})^H \quad (7)$$

$$(\mathbf{A} + \mathbf{B})^H = \mathbf{A}^H + \mathbf{B}^H \quad (8)$$

$$(\mathbf{AB})^H = \mathbf{B}^H\mathbf{A}^H \quad (9)$$

$$(\mathbf{ABC}\dots)^H = \dots\mathbf{C}^H\mathbf{B}^H\mathbf{A}^H \quad (10)$$

1.1 Trace

$$\text{Tr}(\mathbf{A}) = \sum_i A_{ii} \quad (11)$$

$$\text{Tr}(\mathbf{A}) = \sum_i \lambda_i, \quad \lambda_i = \text{eig}(\mathbf{A}) \quad (12)$$

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^T) \quad (13)$$

$$\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA}) \quad (14)$$

$$\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) \quad (15)$$

$$\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{BCA}) = \text{Tr}(\mathbf{CAB}) \quad (16)$$

$$\mathbf{a}^T \mathbf{a} = \text{Tr}(\mathbf{a}\mathbf{a}^T) \quad (17)$$

1.2 Determinant

Let \mathbf{A} be an $n \times n$ matrix.

$$\det(\mathbf{A}) = \prod_i \lambda_i \quad \lambda_i = \text{eig}(\mathbf{A}) \quad (18)$$

$$\det(c\mathbf{A}) = c^n \det(\mathbf{A}), \quad \text{if } \mathbf{A} \in \mathbb{R}^{n \times n} \quad (19)$$

$$\det(\mathbf{A}^T) = \det(\mathbf{A}) \quad (20)$$

$$\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B}) \quad (21)$$

$$\det(\mathbf{A}^{-1}) = 1 / \det(\mathbf{A}) \quad (22)$$

$$\det(\mathbf{A}^n) = \det(\mathbf{A})^n \quad (23)$$

$$\det(\mathbf{I} + \mathbf{u}\mathbf{v}^T) = 1 + \mathbf{u}^T \mathbf{v} \quad (24)$$

For $n = 2$:

$$\det(\mathbf{I} + \mathbf{A}) = 1 + \det(\mathbf{A}) + \text{Tr}(\mathbf{A}) \quad (25)$$

For $n = 3$:

$$\det(\mathbf{I} + \mathbf{A}) = 1 + \det(\mathbf{A}) + \text{Tr}(\mathbf{A}) + \frac{1}{2}\text{Tr}(\mathbf{A})^2 - \frac{1}{2}\text{Tr}(\mathbf{A}^2) \quad (26)$$

For $n = 4$:

$$\begin{aligned}\det(\mathbf{I} + \mathbf{A}) &= 1 + \det(\mathbf{A}) + \text{Tr}(\mathbf{A}) + \frac{1}{2} \\ &\quad + \text{Tr}(\mathbf{A})^2 - \frac{1}{2}\text{Tr}(\mathbf{A}^2) \\ &\quad + \frac{1}{6}\text{Tr}(\mathbf{A})^3 - \frac{1}{2}\text{Tr}(\mathbf{A})\text{Tr}(\mathbf{A}^2) + \frac{1}{3}\text{Tr}(\mathbf{A}^3)\end{aligned}\quad (27)$$

For small ε , the following approximation holds

$$\det(\mathbf{I} + \varepsilon\mathbf{A}) \cong 1 + \det(\mathbf{A}) + \varepsilon\text{Tr}(\mathbf{A}) + \frac{1}{2}\varepsilon^2\text{Tr}(\mathbf{A})^2 - \frac{1}{2}\varepsilon^2\text{Tr}(\mathbf{A}^2) \quad (28)$$

1.3 The Special Case 2x2

Consider the matrix \mathbf{A}

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

Determinant and trace

$$\det(\mathbf{A}) = A_{11}A_{22} - A_{12}A_{21} \quad (29)$$

$$\text{Tr}(\mathbf{A}) = A_{11} + A_{22} \quad (30)$$

Eigenvalues

$$\lambda^2 - \lambda \cdot \text{Tr}(\mathbf{A}) + \det(\mathbf{A}) = 0$$

$$\lambda_1 = \frac{\text{Tr}(\mathbf{A}) + \sqrt{\text{Tr}(\mathbf{A})^2 - 4\det(\mathbf{A})}}{2} \quad \lambda_2 = \frac{\text{Tr}(\mathbf{A}) - \sqrt{\text{Tr}(\mathbf{A})^2 - 4\det(\mathbf{A})}}{2}$$

$$\lambda_1 + \lambda_2 = \text{Tr}(\mathbf{A}) \quad \lambda_1\lambda_2 = \det(\mathbf{A})$$

Eigenvectors

$$\mathbf{v}_1 \propto \begin{bmatrix} A_{12} \\ \lambda_1 - A_{11} \end{bmatrix} \quad \mathbf{v}_2 \propto \begin{bmatrix} A_{12} \\ \lambda_2 - A_{11} \end{bmatrix}$$

Inverse

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \begin{bmatrix} A_{22} & -A_{12} \\ -A_{21} & A_{11} \end{bmatrix} \quad (31)$$

2 Derivatives

This section is covering differentiation of a number of expressions with respect to a matrix \mathbf{X} . Note that it is always assumed that \mathbf{X} has *no special structure*, i.e. that the elements of \mathbf{X} are independent (e.g. not symmetric, Toeplitz, positive definite). See section 2.8 for differentiation of structured matrices. The basic assumptions can be written in a formula as

$$\frac{\partial X_{kl}}{\partial X_{ij}} = \delta_{ik}\delta_{lj} \quad (32)$$

that is for e.g. vector forms,

$$\left[\frac{\partial \mathbf{x}}{\partial y}\right]_i = \frac{\partial x_i}{\partial y} \quad \left[\frac{\partial x}{\partial \mathbf{y}}\right]_i = \frac{\partial x}{\partial y_i} \quad \left[\frac{\partial \mathbf{x}}{\partial \mathbf{y}}\right]_{ij} = \frac{\partial x_i}{\partial y_j}$$

The following rules are general and very useful when deriving the differential of an expression ([19]):

$$\partial \mathbf{A} = 0 \quad (\mathbf{A} \text{ is a constant}) \quad (33)$$

$$\partial(\alpha \mathbf{X}) = \alpha \partial \mathbf{X} \quad (34)$$

$$\partial(\mathbf{X} + \mathbf{Y}) = \partial \mathbf{X} + \partial \mathbf{Y} \quad (35)$$

$$\partial(\text{Tr}(\mathbf{X})) = \text{Tr}(\partial \mathbf{X}) \quad (36)$$

$$\partial(\mathbf{X}\mathbf{Y}) = (\partial \mathbf{X})\mathbf{Y} + \mathbf{X}(\partial \mathbf{Y}) \quad (37)$$

$$\partial(\mathbf{X} \circ \mathbf{Y}) = (\partial \mathbf{X}) \circ \mathbf{Y} + \mathbf{X} \circ (\partial \mathbf{Y}) \quad (38)$$

$$\partial(\mathbf{X} \otimes \mathbf{Y}) = (\partial \mathbf{X}) \otimes \mathbf{Y} + \mathbf{X} \otimes (\partial \mathbf{Y}) \quad (39)$$

$$\partial(\mathbf{X}^{-1}) = -\mathbf{X}^{-1}(\partial \mathbf{X})\mathbf{X}^{-1} \quad (40)$$

$$\partial(\det(\mathbf{X})) = \text{Tr}(\text{adj}(\mathbf{X})\partial \mathbf{X}) \quad (41)$$

$$\partial(\det(\mathbf{X})) = \det(\mathbf{X})\text{Tr}(\mathbf{X}^{-1}\partial \mathbf{X}) \quad (42)$$

$$\partial(\ln(\det(\mathbf{X}))) = \text{Tr}(\mathbf{X}^{-1}\partial \mathbf{X}) \quad (43)$$

$$\partial \mathbf{X}^T = (\partial \mathbf{X})^T \quad (44)$$

$$\partial \mathbf{X}^H = (\partial \mathbf{X})^H \quad (45)$$

2.1 Derivatives of a Determinant

2.1.1 General form

$$\frac{\partial \det(\mathbf{Y})}{\partial x} = \det(\mathbf{Y})\text{Tr}\left[\mathbf{Y}^{-1}\frac{\partial \mathbf{Y}}{\partial x}\right] \quad (46)$$

$$\sum_k \frac{\partial \det(\mathbf{X})}{\partial X_{ik}} X_{jk} = \delta_{ij} \det(\mathbf{X}) \quad (47)$$

$$\begin{aligned} \frac{\partial^2 \det(\mathbf{Y})}{\partial x^2} &= \det(\mathbf{Y}) \left[\text{Tr} \left[\mathbf{Y}^{-1} \frac{\partial^2 \mathbf{Y}}{\partial x^2} \right] \right. \\ &\quad + \text{Tr} \left[\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \right] \text{Tr} \left[\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \right] \\ &\quad \left. - \text{Tr} \left[\left(\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \right) \left(\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \right) \right] \right] \quad (48) \end{aligned}$$

2.1.2 Linear forms

$$\frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \det(\mathbf{X})(\mathbf{X}^{-1})^T \quad (49)$$

$$\sum_k \frac{\partial \det(\mathbf{X})}{\partial X_{ik}} X_{jk} = \delta_{ij} \det(\mathbf{X}) \quad (50)$$

$$\frac{\partial \det(\mathbf{AXB})}{\partial \mathbf{X}} = \det(\mathbf{AXB})(\mathbf{X}^{-1})^T = \det(\mathbf{AXB})(\mathbf{X}^T)^{-1} \quad (51)$$

2.1.3 Square forms

If \mathbf{X} is square and invertible, then

$$\frac{\partial \det(\mathbf{X}^T \mathbf{A} \mathbf{X})}{\partial \mathbf{X}} = 2 \det(\mathbf{X}^T \mathbf{A} \mathbf{X}) \mathbf{X}^{-T} \quad (52)$$

If \mathbf{X} is not square but \mathbf{A} is symmetric, then

$$\frac{\partial \det(\mathbf{X}^T \mathbf{A} \mathbf{X})}{\partial \mathbf{X}} = 2 \det(\mathbf{X}^T \mathbf{A} \mathbf{X}) \mathbf{A} \mathbf{X} (\mathbf{X}^T \mathbf{A} \mathbf{X})^{-1} \quad (53)$$

If \mathbf{X} is not square and \mathbf{A} is not symmetric, then

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

2.1.4 Other nonlinear forms

Some special cases are (See [9, 7])

$$\frac{\partial \ln \det(\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}} = 2(\mathbf{X}^+)^T \quad (55)$$

$$\frac{\partial \ln \det(\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}^+} = -2\mathbf{X}^T \quad (56)$$

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

$$\frac{\partial \det(\mathbf{X}^k)}{\partial \mathbf{X}} = k \det(\mathbf{X}^k) \mathbf{X}^{-T} \quad (58)$$

2.2 Derivatives of an Inverse

From [27] we have the basic identity

$$\frac{\partial \mathbf{Y}^{-1}}{\partial x} = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \mathbf{Y}^{-1} \quad (59)$$

from which it follows

$$\frac{\partial(\mathbf{X}^{-1})_{kl}}{\partial X_{ij}} = -(\mathbf{X}^{-1})_{ki}(\mathbf{X}^{-1})_{jl} \quad (60)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T} \quad (61)$$

$$\frac{\partial \det(\mathbf{X}^{-1})}{\partial \mathbf{X}} = -\det(\mathbf{X}^{-1})(\mathbf{X}^{-1})^T \quad (62)$$

$$\frac{\partial \text{Tr}(\mathbf{A} \mathbf{X}^{-1} \mathbf{B})}{\partial \mathbf{X}} = -(\mathbf{X}^{-1} \mathbf{B} \mathbf{A} \mathbf{X}^{-1})^T \quad (63)$$

$$\frac{\partial \text{Tr}((\mathbf{X} + \mathbf{A})^{-1})}{\partial \mathbf{X}} = -((\mathbf{X} + \mathbf{A})^{-1}(\mathbf{X} + \mathbf{A})^{-1})^T \quad (64)$$

From [32] we have the following result: Let \mathbf{A} be an $n \times n$ invertible square matrix, \mathbf{W} be the inverse of \mathbf{A} , and $J(\mathbf{A})$ is an $n \times n$ -variate and differentiable function with respect to \mathbf{A} , then the partial differentials of J with respect to \mathbf{A} and \mathbf{W} satisfy

$$\frac{\partial J}{\partial \mathbf{A}} = -\mathbf{A}^{-T} \frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T}$$

2.3 Derivatives of Eigenvalues

$$\frac{\partial}{\partial \mathbf{X}} \sum \text{eig}(\mathbf{X}) = \frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}) = \mathbf{I} \quad (65)$$

$$\frac{\partial}{\partial \mathbf{X}} \prod \text{eig}(\mathbf{X}) = \frac{\partial}{\partial \mathbf{X}} \det(\mathbf{X}) = \det(\mathbf{X}) \mathbf{X}^{-T} \quad (66)$$

If \mathbf{A} is real and symmetric, λ_i and \mathbf{v}_i are distinct eigenvalues and eigenvectors of \mathbf{A} (see (276)) with $\mathbf{v}_i^T \mathbf{v}_i = 1$, then [33]

$$\partial \lambda_i = \mathbf{v}_i^T \partial(\mathbf{A}) \mathbf{v}_i \quad (67)$$

$$\partial \mathbf{v}_i = (\lambda_i \mathbf{I} - \mathbf{A})^+ \partial(\mathbf{A}) \mathbf{v}_i \quad (68)$$

2.4 Derivatives of Matrices, Vectors and Scalar Forms

2.4.1 First Order

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad (69)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T \quad (70)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^T \quad (71)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} = \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^T \quad (72)$$

$$\frac{\partial \mathbf{X}}{\partial X_{ij}} = \mathbf{J}^{ij} \quad (73)$$

$$\frac{\partial (\mathbf{X} \mathbf{A})_{ij}}{\partial X_{mn}} = \delta_{im} (\mathbf{A})_{nj} = (\mathbf{J}^{mn} \mathbf{A})_{ij} \quad (74)$$

$$\frac{\partial (\mathbf{X}^T \mathbf{A})_{ij}}{\partial X_{mn}} = \delta_{in} (\mathbf{A})_{mj} = (\mathbf{J}^{nm} \mathbf{A})_{ij} \quad (75)$$

2.4.2 Second Order

$$\frac{\partial}{\partial X_{ij}} \sum_{klmn} X_{kl} X_{mn} = 2 \sum_{kl} X_{kl} \quad (76)$$

$$\frac{\partial \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{c}}{\partial \mathbf{X}} = \mathbf{X}(\mathbf{b} \mathbf{c}^T + \mathbf{c} \mathbf{b}^T) \quad (77)$$

$$\frac{\partial (\mathbf{B} \mathbf{x} + \mathbf{b})^T \mathbf{C} (\mathbf{D} \mathbf{x} + \mathbf{d})}{\partial \mathbf{x}} = \mathbf{B}^T \mathbf{C} (\mathbf{D} \mathbf{x} + \mathbf{d}) + \mathbf{D}^T \mathbf{C}^T (\mathbf{B} \mathbf{x} + \mathbf{b}) \quad (78)$$

$$\frac{\partial (\mathbf{X}^T \mathbf{B} \mathbf{X})_{kl}}{\partial X_{ij}} = \delta_{lj} (\mathbf{X}^T \mathbf{B})_{ki} + \delta_{ki} (\mathbf{B} \mathbf{X})_{jl} \quad (79)$$

$$\frac{\partial (\mathbf{X}^T \mathbf{B} \mathbf{X})}{\partial X_{ij}} = \mathbf{X}^T \mathbf{B} \mathbf{J}^{ij} + \mathbf{J}^{ji} \mathbf{B} \mathbf{X} \quad (\mathbf{J}^{ij})_{kl} = \delta_{ik} \delta_{jl} \quad (80)$$

See Sec 9.7 for useful properties of the Single-entry matrix \mathbf{J}^{ij}

$$\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x} \quad (81)$$

$$\frac{\partial \mathbf{b}^T \mathbf{X}^T \mathbf{D} \mathbf{X} \mathbf{c}}{\partial \mathbf{X}} = \mathbf{D}^T \mathbf{X} \mathbf{b} \mathbf{c}^T + \mathbf{D} \mathbf{X} \mathbf{c} \mathbf{b}^T \quad (82)$$

$$\frac{\partial}{\partial \mathbf{X}} (\mathbf{X} \mathbf{b} + \mathbf{c})^T \mathbf{D} (\mathbf{X} \mathbf{b} + \mathbf{c}) = (\mathbf{D} + \mathbf{D}^T) (\mathbf{X} \mathbf{b} + \mathbf{c}) \mathbf{b}^T \quad (83)$$

Assume \mathbf{W} is symmetric, then

$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = -2 \mathbf{A}^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \quad (84)$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} - \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{s}) = 2 \mathbf{W} (\mathbf{x} - \mathbf{s}) \quad (85)$$

$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{s}) = -2 \mathbf{W} (\mathbf{x} - \mathbf{s}) \quad (86)$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = 2 \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \quad (87)$$

$$\frac{\partial}{\partial \mathbf{A}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = -2 \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \mathbf{s}^T \quad (88)$$

As a case with complex values the following holds

$$\frac{\partial (a - \mathbf{x}^H \mathbf{b})^2}{\partial \mathbf{x}} = -2 \mathbf{b} (a - \mathbf{x}^H \mathbf{b})^* \quad (89)$$

This formula is also known from the LMS algorithm [14]

2.4.3 Higher-order and non-linear

$$\frac{\partial (\mathbf{X}^n)_{kl}}{\partial X_{ij}} = \sum_{r=0}^{n-1} (\mathbf{X}^r \mathbf{J}^{ij} \mathbf{X}^{n-1-r})_{kl} \quad (90)$$

For proof of the above, see B.1.3.

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{a}^T \mathbf{X}^n \mathbf{b} = \sum_{r=0}^{n-1} (\mathbf{X}^r)^T \mathbf{a} \mathbf{b}^T (\mathbf{X}^{n-1-r})^T \quad (91)$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{X}} \mathbf{a}^T (\mathbf{X}^n)^T \mathbf{X}^n \mathbf{b} &= \sum_{r=0}^{n-1} \left[\mathbf{X}^{n-1-r} \mathbf{a} \mathbf{b}^T (\mathbf{X}^n)^T \mathbf{X}^r \right. \\ &\quad \left. + (\mathbf{X}^r)^T \mathbf{X}^n \mathbf{a} \mathbf{b}^T (\mathbf{X}^{n-1-r})^T \right] \end{aligned} \quad (92)$$

See B.1.3 for a proof.

Assume \mathbf{s} and \mathbf{r} are functions of \mathbf{x} , i.e. $\mathbf{s} = \mathbf{s}(\mathbf{x})$, $\mathbf{r} = \mathbf{r}(\mathbf{x})$, and that \mathbf{A} is a constant, then

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{s}^T \mathbf{A} \mathbf{r} = \left[\frac{\partial \mathbf{s}}{\partial \mathbf{x}} \right]^T \mathbf{A} \mathbf{r} + \left[\frac{\partial \mathbf{r}}{\partial \mathbf{x}} \right]^T \mathbf{A}^T \mathbf{s} \quad (93)$$

$$\frac{\partial}{\partial \mathbf{x}} \frac{(\mathbf{A} \mathbf{x})^T (\mathbf{A} \mathbf{x})}{(\mathbf{B} \mathbf{x})^T (\mathbf{B} \mathbf{x})} = \frac{\partial}{\partial \mathbf{x}} \frac{\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B}^T \mathbf{B} \mathbf{x}} \quad (94)$$

$$= 2 \frac{\mathbf{A}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{B} \mathbf{x}} - 2 \frac{\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} \mathbf{B}^T \mathbf{B} \mathbf{x}}{(\mathbf{x}^T \mathbf{B}^T \mathbf{B} \mathbf{x})^2} \quad (95)$$

2.4.4 Gradient and Hessian

Using the above we have for the gradient and the Hessian

$$f = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} \quad (96)$$

$$\nabla_{\mathbf{x}} f = \frac{\partial f}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} + \mathbf{b} \quad (97)$$

$$\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^T} = \mathbf{A} + \mathbf{A}^T \quad (98)$$

2.5 Derivatives of Traces

Assume $F(\mathbf{X})$ to be a differentiable function of each of the elements of \mathbf{X} . It then holds that

$$\frac{\partial \text{Tr}(F(\mathbf{X}))}{\partial \mathbf{X}} = f(\mathbf{X})^T$$

where $f(\cdot)$ is the scalar derivative of $F(\cdot)$.

2.5.1 First Order

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}) = \mathbf{I} \quad (99)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X} \mathbf{A}) = \mathbf{A}^T \quad (100)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A} \mathbf{X} \mathbf{B}) = \mathbf{A}^T \mathbf{B}^T \quad (101)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A} \mathbf{X}^T \mathbf{B}) = \mathbf{B} \mathbf{A} \quad (102)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{A}) = \mathbf{A} \quad (103)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A} \mathbf{X}^T) = \mathbf{A} \quad (104)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A} \otimes \mathbf{X}) = \text{Tr}(\mathbf{A}) \mathbf{I} \quad (105)$$

2.5.2 Second Order

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^2) = 2\mathbf{X}^T \quad (106)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^2 \mathbf{B}) = (\mathbf{X}\mathbf{B} + \mathbf{B}\mathbf{X})^T \quad (107)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{B}\mathbf{X}) = \mathbf{B}\mathbf{X} + \mathbf{B}^T \mathbf{X} \quad (108)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{B}\mathbf{X}\mathbf{X}^T) = \mathbf{B}\mathbf{X} + \mathbf{B}^T \mathbf{X} \quad (109)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{X}^T \mathbf{B}) = \mathbf{B}\mathbf{X} + \mathbf{B}^T \mathbf{X} \quad (110)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{B}\mathbf{X}^T) = \mathbf{X}\mathbf{B}^T + \mathbf{X}\mathbf{B} \quad (111)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{B}\mathbf{X}^T \mathbf{X}) = \mathbf{X}\mathbf{B}^T + \mathbf{X}\mathbf{B} \quad (112)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{X}\mathbf{B}) = \mathbf{X}\mathbf{B}^T + \mathbf{X}\mathbf{B} \quad (113)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}) = \mathbf{A}^T \mathbf{X}^T \mathbf{B}^T + \mathbf{B}^T \mathbf{X}^T \mathbf{A}^T \quad (114)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{X}) = \frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{X}^T) = 2\mathbf{X} \quad (115)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{B}^T \mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{B}) = \mathbf{C}^T \mathbf{X}\mathbf{B}\mathbf{B}^T + \mathbf{C}\mathbf{X}\mathbf{B}\mathbf{B}^T \quad (116)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}[\mathbf{X}^T \mathbf{B}\mathbf{X}\mathbf{C}] = \mathbf{B}\mathbf{X}\mathbf{C} + \mathbf{B}^T \mathbf{X}\mathbf{C}^T \quad (117)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}^T \mathbf{C}) = \mathbf{A}^T \mathbf{C}^T \mathbf{X}\mathbf{B}^T + \mathbf{C}\mathbf{A}\mathbf{X}\mathbf{B} \quad (118)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}[(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})^T] = 2\mathbf{A}^T (\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})\mathbf{B}^T \quad (119)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X} \otimes \mathbf{X}) = \frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X})\text{Tr}(\mathbf{X}) = 2\text{Tr}(\mathbf{X})\mathbf{I} \quad (120)$$

See [7].

2.5.3 Higher Order

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^k) = k(\mathbf{X}^{k-1})^T \quad (121)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A}\mathbf{X}^k) = \sum_{r=0}^{k-1} (\mathbf{X}^r \mathbf{A} \mathbf{X}^{k-r-1})^T \quad (122)$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{X}} \text{Tr}[\mathbf{B}^T \mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{B}] &= \mathbf{C}\mathbf{X}\mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{B}\mathbf{B}^T \\ &\quad + \mathbf{C}^T \mathbf{X}\mathbf{B}\mathbf{B}^T \mathbf{X}^T \mathbf{C}^T \mathbf{X} \\ &\quad + \mathbf{C}\mathbf{X}\mathbf{B}\mathbf{B}^T \mathbf{X}^T \mathbf{C}\mathbf{X} \\ &\quad + \mathbf{C}^T \mathbf{X}\mathbf{X}^T \mathbf{C}^T \mathbf{X}\mathbf{B}\mathbf{B}^T \end{aligned} \quad (123)$$

2.5.4 Other

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A}\mathbf{X}^{-1}\mathbf{B}) = -(\mathbf{X}^{-1}\mathbf{B}\mathbf{A}\mathbf{X}^{-1})^T = -\mathbf{X}^{-T}\mathbf{A}^T\mathbf{B}^T\mathbf{X}^{-T} \quad (124)$$

Assume \mathbf{B} and \mathbf{C} to be symmetric, then

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}[(\mathbf{X}^T\mathbf{C}\mathbf{X})^{-1}\mathbf{A}] = -(\mathbf{C}\mathbf{X}(\mathbf{X}^T\mathbf{C}\mathbf{X})^{-1})(\mathbf{A} + \mathbf{A}^T)(\mathbf{X}^T\mathbf{C}\mathbf{X})^{-1} \quad (125)$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{X}} \text{Tr}[(\mathbf{X}^T\mathbf{C}\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{B}\mathbf{X})] &= -2\mathbf{C}\mathbf{X}(\mathbf{X}^T\mathbf{C}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{B}\mathbf{X}(\mathbf{X}^T\mathbf{C}\mathbf{X})^{-1} \\ &\quad + 2\mathbf{B}\mathbf{X}(\mathbf{X}^T\mathbf{C}\mathbf{X})^{-1} \end{aligned} \quad (126)$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{X}} \text{Tr}[(\mathbf{A} + \mathbf{X}^T\mathbf{C}\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{B}\mathbf{X})] &= -2\mathbf{C}\mathbf{X}(\mathbf{A} + \mathbf{X}^T\mathbf{C}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{B}\mathbf{X}(\mathbf{A} + \mathbf{X}^T\mathbf{C}\mathbf{X})^{-1} \\ &\quad + 2\mathbf{B}\mathbf{X}(\mathbf{A} + \mathbf{X}^T\mathbf{C}\mathbf{X})^{-1} \end{aligned} \quad (127)$$

See [7].

$$\frac{\partial \text{Tr}(\sin(\mathbf{X}))}{\partial \mathbf{X}} = \cos(\mathbf{X})^T \quad (128)$$

2.6 Derivatives of vector norms**2.6.1 Two-norm**

$$\frac{\partial}{\partial \mathbf{x}} \|\mathbf{x} - \mathbf{a}\|_2 = \frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|_2} \quad (129)$$

$$\frac{\partial}{\partial \mathbf{x}} \frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|_2} = \frac{\mathbf{I}}{\|\mathbf{x} - \mathbf{a}\|_2} - \frac{(\mathbf{x} - \mathbf{a})(\mathbf{x} - \mathbf{a})^T}{\|\mathbf{x} - \mathbf{a}\|_2^3} \quad (130)$$

$$\frac{\partial \|\mathbf{x}\|_2^2}{\partial \mathbf{x}} = \frac{\partial \|\mathbf{x}^T \mathbf{x}\|_2}{\partial \mathbf{x}} = 2\mathbf{x} \quad (131)$$

2.7 Derivatives of matrix norms

For more on matrix norms, see Sec. 10.4.

2.7.1 Frobenius norm

$$\frac{\partial}{\partial \mathbf{X}} \|\mathbf{X}\|_F^2 = \frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{X}^H) = 2\mathbf{X} \quad (132)$$

See (248). Note that this is also a special case of the result in equation 119.

2.8 Derivatives of Structured Matrices

Assume that the matrix \mathbf{A} has some structure, i.e. symmetric, toeplitz, etc. In that case the derivatives of the previous section does not apply in general. Instead, consider the following general rule for differentiating a scalar function $f(\mathbf{A})$

$$\frac{df}{dA_{ij}} = \sum_{kl} \frac{\partial f}{\partial A_{kl}} \frac{\partial A_{kl}}{\partial A_{ij}} = \text{Tr} \left[\left[\frac{\partial f}{\partial \mathbf{A}} \right]^T \frac{\partial \mathbf{A}}{\partial A_{ij}} \right] \quad (133)$$

The matrix differentiated with respect to itself is in this document referred to as the *structure matrix* of \mathbf{A} and is defined simply by

$$\frac{\partial \mathbf{A}}{\partial A_{ij}} = \mathbf{S}^{ij} \quad (134)$$

If \mathbf{A} has no special structure we have simply $\mathbf{S}^{ij} = \mathbf{J}^{ij}$, that is, the structure matrix is simply the single-entry matrix. Many structures have a representation in singleentry matrices, see Sec. 9.7.6 for more examples of structure matrices.

2.8.1 The Chain Rule

Sometimes the objective is to find the derivative of a matrix which is a function of another matrix. Let $\mathbf{U} = f(\mathbf{X})$, the goal is to find the derivative of the function $g(\mathbf{U})$ with respect to \mathbf{X} :

$$\frac{\partial g(\mathbf{U})}{\partial \mathbf{X}} = \frac{\partial g(f(\mathbf{X}))}{\partial \mathbf{X}} \quad (135)$$

Then the Chain Rule can then be written the following way:

$$\frac{\partial g(\mathbf{U})}{\partial \mathbf{X}} = \frac{\partial g(\mathbf{U})}{\partial x_{ij}} = \sum_{k=1}^M \sum_{l=1}^N \frac{\partial g(\mathbf{U})}{\partial u_{kl}} \frac{\partial u_{kl}}{\partial x_{ij}} \quad (136)$$

Using matrix notation, this can be written as:

$$\frac{\partial g(\mathbf{U})}{\partial X_{ij}} = \text{Tr} \left[\left(\frac{\partial g(\mathbf{U})}{\partial \mathbf{U}} \right)^T \frac{\partial \mathbf{U}}{\partial X_{ij}} \right]. \quad (137)$$

2.8.2 Symmetric

If \mathbf{A} is symmetric, then $\mathbf{S}^{ij} = \mathbf{J}^{ij} + \mathbf{J}^{ji} - \mathbf{J}^{ij} \mathbf{J}^{ij}$ and therefore

$$\frac{df}{d\mathbf{A}} = \left[\frac{\partial f}{\partial \mathbf{A}} \right] + \left[\frac{\partial f}{\partial \mathbf{A}} \right]^T - \text{diag} \left[\frac{\partial f}{\partial \mathbf{A}} \right] \quad (138)$$

That is, e.g., ([5]):

$$\frac{\partial \text{Tr}(\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \mathbf{A} + \mathbf{A}^T - (\mathbf{A} \circ \mathbf{I}), \text{ see (142)} \quad (139)$$

$$\frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \det(\mathbf{X})(2\mathbf{X}^{-1} - (\mathbf{X}^{-1} \circ \mathbf{I})) \quad (140)$$

$$\frac{\partial \ln \det(\mathbf{X})}{\partial \mathbf{X}} = 2\mathbf{X}^{-1} - (\mathbf{X}^{-1} \circ \mathbf{I}) \quad (141)$$

2.8.3 Diagonal

If \mathbf{X} is diagonal, then ([19]):

$$\frac{\partial \text{Tr}(\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \mathbf{A} \circ \mathbf{I} \quad (142)$$

2.8.4 Toeplitz

Like symmetric matrices and diagonal matrices also Toeplitz matrices has a special structure which should be taken into account when the derivative with respect to a matrix with Toeplitz structure.

$$\begin{aligned}
 & \frac{\partial \text{Tr}(\mathbf{A}\mathbf{T})}{\partial \mathbf{T}} \\
 &= \frac{\partial \text{Tr}(\mathbf{T}\mathbf{A})}{\partial \mathbf{T}} \\
 &= \begin{bmatrix} \text{Tr}(\mathbf{A}) & \text{Tr}([\mathbf{A}^T]_{n1}) & \text{Tr}([\mathbf{A}^T]_{1n}]_{n-1,2}) & \cdots & A_{n1} \\ \text{Tr}([\mathbf{A}^T]_{1n}) & \text{Tr}(\mathbf{A}) & \ddots & \ddots & \vdots \\ \text{Tr}([\mathbf{A}^T]_{1n}]_{2,n-1}) & \ddots & \ddots & \ddots & \text{Tr}([\mathbf{A}^T]_{1n}]_{n-1,2}) \\ \vdots & \ddots & \ddots & \ddots & \text{Tr}([\mathbf{A}^T]_{n1}) \\ A_{1n} & \cdots & \text{Tr}([\mathbf{A}^T]_{1n}]_{2,n-1}) & \text{Tr}([\mathbf{A}^T]_{1n}) & \text{Tr}(\mathbf{A}) \end{bmatrix} \\
 &\equiv \boldsymbol{\alpha}(\mathbf{A})
 \end{aligned} \tag{143}$$

As it can be seen, the derivative $\boldsymbol{\alpha}(\mathbf{A})$ also has a Toeplitz structure. Each value in the diagonal is the sum of all the diagonal valued in \mathbf{A} , the values in the diagonals next to the main diagonal equal the sum of the diagonal next to the main diagonal in \mathbf{A}^T . This result is only valid for the unconstrained Toeplitz matrix. If the Toeplitz matrix also is symmetric, the same derivative yields

$$\frac{\partial \text{Tr}(\mathbf{A}\mathbf{T})}{\partial \mathbf{T}} = \frac{\partial \text{Tr}(\mathbf{T}\mathbf{A})}{\partial \mathbf{T}} = \boldsymbol{\alpha}(\mathbf{A}) + \boldsymbol{\alpha}(\mathbf{A})^T - \boldsymbol{\alpha}(\mathbf{A}) \circ \mathbf{I} \tag{144}$$

3 Inverses

3.1 Basic

3.1.1 Definition

The *inverse* \mathbf{A}^{-1} of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is defined such that

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}, \quad (145)$$

where \mathbf{I} is the $n \times n$ identity matrix. If \mathbf{A}^{-1} exists, \mathbf{A} is said to be *nonsingular*. Otherwise, \mathbf{A} is said to be *singular* (see e.g. [12]).

3.1.2 Cofactors and Adjoint

The *submatrix* of a matrix \mathbf{A} , denoted by $[\mathbf{A}]_{ij}$ is a $(n-1) \times (n-1)$ matrix obtained by deleting the i th row and the j th column of \mathbf{A} . The (i, j) *cofactor* of a matrix is defined as

$$\text{cof}(\mathbf{A}, i, j) = (-1)^{i+j} \det([\mathbf{A}]_{ij}), \quad (146)$$

The *matrix of cofactors* can be created from the cofactors

$$\text{cof}(\mathbf{A}) = \begin{bmatrix} \text{cof}(\mathbf{A}, 1, 1) & \cdots & \text{cof}(\mathbf{A}, 1, n) \\ \vdots & \text{cof}(\mathbf{A}, i, j) & \vdots \\ \text{cof}(\mathbf{A}, n, 1) & \cdots & \text{cof}(\mathbf{A}, n, n) \end{bmatrix} \quad (147)$$

The *adjoint* matrix is the transpose of the cofactor matrix

$$\text{adj}(\mathbf{A}) = (\text{cof}(\mathbf{A}))^T, \quad (148)$$

3.1.3 Determinant

The *determinant* of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is defined as (see [12])

$$\det(\mathbf{A}) = \sum_{j=1}^n (-1)^{j+1} A_{1j} \det([\mathbf{A}]_{1j}) \quad (149)$$

$$= \sum_{j=1}^n A_{1j} \text{cof}(\mathbf{A}, 1, j). \quad (150)$$

3.1.4 Construction

The inverse matrix can be constructed, using the adjoint matrix, by

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \cdot \text{adj}(\mathbf{A}) \quad (151)$$

For the case of 2×2 matrices, see section 1.3.

3.1.5 Condition number

The condition number of a matrix $c(\mathbf{A})$ is the ratio between the largest and the smallest singular value of a matrix (see Section 5.3 on singular values),

$$c(\mathbf{A}) = \frac{d_+}{d_-} \quad (152)$$

The condition number can be used to measure how singular a matrix is. If the condition number is large, it indicates that the matrix is nearly singular. The condition number can also be estimated from the matrix norms. Here

$$c(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|, \quad (153)$$

where $\|\cdot\|$ is a norm such as e.g the 1-norm, the 2-norm, the ∞ -norm or the Frobenius norm (see Sec 10.4 for more on matrix norms).

The 2-norm of \mathbf{A} equals $\sqrt{(\max(\text{eig}(\mathbf{A}^H \mathbf{A}))})}$ [12, p.57]. For a symmetric matrix, this reduces to $\|\mathbf{A}\|_2 = \max(|\text{eig}(\mathbf{A})|)$ [12, p.394]. If the matrix is symmetric and positive definite, $\|\mathbf{A}\|_2 = \max(\text{eig}(\mathbf{A}))$. The condition number based on the 2-norm thus reduces to

$$\|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \max(\text{eig}(\mathbf{A})) \max(\text{eig}(\mathbf{A}^{-1})) = \frac{\max(\text{eig}(\mathbf{A}))}{\min(\text{eig}(\mathbf{A}))}. \quad (154)$$

3.2 Exact Relations

3.2.1 Basic

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1} \quad (155)$$

3.2.2 The Woodbury identity

The Woodbury identity comes in many variants. The latter of the two can be found in [12]

$$(\mathbf{A} + \mathbf{CBC}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{C} (\mathbf{B}^{-1} + \mathbf{C}^T \mathbf{A}^{-1} \mathbf{C})^{-1} \mathbf{C}^T \mathbf{A}^{-1} \quad (156)$$

$$(\mathbf{A} + \mathbf{UBV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{B}^{-1} + \mathbf{VA}^{-1} \mathbf{U})^{-1} \mathbf{VA}^{-1} \quad (157)$$

If \mathbf{P}, \mathbf{R} are positive definite, then (see [30])

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} = \mathbf{PB}^T (\mathbf{BPB}^T + \mathbf{R})^{-1} \quad (158)$$

3.2.3 The Kailath Variant

$$(\mathbf{A} + \mathbf{BC})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{I} + \mathbf{CA}^{-1} \mathbf{B})^{-1} \mathbf{CA}^{-1} \quad (159)$$

See [4, page 153].

3.2.4 Sherman-Morrison

$$(\mathbf{A} + \mathbf{bc}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{bc}^T \mathbf{A}^{-1}}{1 + \mathbf{c}^T \mathbf{A}^{-1} \mathbf{b}} \quad (160)$$

3.2.5 The Searle Set of Identities

The following set of identities, can be found in [25, page 151],

$$(\mathbf{I} + \mathbf{A}^{-1})^{-1} = \mathbf{A}(\mathbf{A} + \mathbf{I})^{-1} \quad (161)$$

$$(\mathbf{A} + \mathbf{B}\mathbf{B}^T)^{-1}\mathbf{B} = \mathbf{A}^{-1}\mathbf{B}(\mathbf{I} + \mathbf{B}^T\mathbf{A}^{-1}\mathbf{B})^{-1} \quad (162)$$

$$(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} = \mathbf{A}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{B} = \mathbf{B}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{A} \quad (163)$$

$$\mathbf{A} - \mathbf{A}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{A} = \mathbf{B} - \mathbf{B}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{B} \quad (164)$$

$$\mathbf{A}^{-1} + \mathbf{B}^{-1} = \mathbf{A}^{-1}(\mathbf{A} + \mathbf{B})\mathbf{B}^{-1} \quad (165)$$

$$(\mathbf{I} + \mathbf{A}\mathbf{B})^{-1} = \mathbf{I} - \mathbf{A}(\mathbf{I} + \mathbf{B}\mathbf{A})^{-1}\mathbf{B} \quad (166)$$

$$(\mathbf{I} + \mathbf{A}\mathbf{B})^{-1}\mathbf{A} = \mathbf{A}(\mathbf{I} + \mathbf{B}\mathbf{A})^{-1} \quad (167)$$

3.2.6 Rank-1 update of inverse of inner product

Denote $\mathbf{A} = (\mathbf{X}^T\mathbf{X})^{-1}$ and that \mathbf{X} is extended to include a new column vector in the end $\tilde{\mathbf{X}} = [\mathbf{X} \ \mathbf{v}]$. Then [34]

$$(\tilde{\mathbf{X}}^T\tilde{\mathbf{X}})^{-1} = \begin{bmatrix} \mathbf{A} + \frac{\mathbf{A}\mathbf{X}^T\mathbf{v}\mathbf{v}^T\mathbf{X}\mathbf{A}^T}{\mathbf{v}^T\mathbf{v} - \mathbf{v}^T\mathbf{X}\mathbf{A}\mathbf{X}^T\mathbf{v}} & \frac{-\mathbf{A}\mathbf{X}^T\mathbf{v}}{\mathbf{v}^T\mathbf{v} - \mathbf{v}^T\mathbf{X}\mathbf{A}\mathbf{X}^T\mathbf{v}} \\ \frac{-\mathbf{v}^T\mathbf{X}\mathbf{A}^T}{\mathbf{v}^T\mathbf{v} - \mathbf{v}^T\mathbf{X}\mathbf{A}\mathbf{X}^T\mathbf{v}} & \frac{1}{\mathbf{v}^T\mathbf{v} - \mathbf{v}^T\mathbf{X}\mathbf{A}\mathbf{X}^T\mathbf{v}} \end{bmatrix}$$

3.2.7 Rank-1 update of Moore-Penrose Inverse

The following is a rank-1 update for the Moore-Penrose pseudo-inverse of real valued matrices and proof can be found in [18]. The matrix \mathbf{G} is defined below:

$$(\mathbf{A} + \mathbf{c}\mathbf{d}^T)^+ = \mathbf{A}^+ + \mathbf{G} \quad (168)$$

Using the the notation

$$\beta = 1 + \mathbf{d}^T\mathbf{A}^+\mathbf{c} \quad (169)$$

$$\mathbf{v} = \mathbf{A}^+\mathbf{c} \quad (170)$$

$$\mathbf{n} = (\mathbf{A}^+)^T\mathbf{d} \quad (171)$$

$$\mathbf{w} = (\mathbf{I} - \mathbf{A}\mathbf{A}^+)\mathbf{c} \quad (172)$$

$$\mathbf{m} = (\mathbf{I} - \mathbf{A}^+\mathbf{A})^T\mathbf{d} \quad (173)$$

the solution is given as six different cases, depending on the entities $\|\mathbf{w}\|$, $\|\mathbf{m}\|$, and β . Please note, that for any (column) vector \mathbf{v} it holds that $\mathbf{v}^+ = \mathbf{v}^T(\mathbf{v}^T\mathbf{v})^{-1} = \frac{\mathbf{v}^T}{\|\mathbf{v}\|^2}$. The solution is:

Case 1 of 6: If $\|\mathbf{w}\| \neq 0$ and $\|\mathbf{m}\| \neq 0$. Then

$$\mathbf{G} = -\mathbf{v}\mathbf{w}^+ - (\mathbf{m}^+)^T\mathbf{n}^T + \beta(\mathbf{m}^+)^T\mathbf{w}^+ \quad (174)$$

$$= -\frac{1}{\|\mathbf{w}\|^2}\mathbf{v}\mathbf{w}^T - \frac{1}{\|\mathbf{m}\|^2}\mathbf{m}\mathbf{n}^T + \frac{\beta}{\|\mathbf{m}\|^2\|\mathbf{w}\|^2}\mathbf{m}\mathbf{w}^T \quad (175)$$

Case 2 of 6: If $\|\mathbf{w}\| = 0$ and $\|\mathbf{m}\| \neq 0$ and $\beta = 0$. Then

$$\mathbf{G} = -\mathbf{v}\mathbf{v}^+\mathbf{A}^+ - (\mathbf{m}^+)^T\mathbf{n}^T \quad (176)$$

$$= -\frac{1}{\|\mathbf{v}\|^2}\mathbf{v}\mathbf{v}^T\mathbf{A}^+ - \frac{1}{\|\mathbf{m}\|^2}\mathbf{m}\mathbf{n}^T \quad (177)$$

Case 3 of 6: If $\|\mathbf{w}\| = 0$ and $\beta \neq 0$. Then

$$\mathbf{G} = \frac{1}{\beta} \mathbf{m} \mathbf{v}^T \mathbf{A}^+ - \frac{\beta}{\|\mathbf{v}\|^2 \|\mathbf{m}\|^2 + |\beta|^2} \left(\frac{\|\mathbf{v}\|^2}{\beta} \mathbf{m} + \mathbf{v} \right) \left(\frac{\|\mathbf{m}\|^2}{\beta} (\mathbf{A}^+)^T \mathbf{v} + \mathbf{n} \right)^T \quad (178)$$

Case 4 of 6: If $\|\mathbf{w}\| \neq 0$ and $\|\mathbf{m}\| = 0$ and $\beta = 0$. Then

$$\mathbf{G} = -\mathbf{A}^+ \mathbf{n} \mathbf{n}^+ - \mathbf{v} \mathbf{w}^+ \quad (179)$$

$$= -\frac{1}{\|\mathbf{n}\|^2} \mathbf{A}^+ \mathbf{n} \mathbf{n}^T - \frac{1}{\|\mathbf{w}\|^2} \mathbf{v} \mathbf{w}^T \quad (180)$$

Case 5 of 6: If $\|\mathbf{m}\| = 0$ and $\beta \neq 0$. Then

$$\mathbf{G} = \frac{1}{\beta} \mathbf{A}^+ \mathbf{n} \mathbf{w}^T - \frac{\beta}{\|\mathbf{n}\|^2 \|\mathbf{w}\|^2 + |\beta|^2} \left(\frac{\|\mathbf{w}\|^2}{\beta} \mathbf{A}^+ \mathbf{n} + \mathbf{v} \right) \left(\frac{\|\mathbf{n}\|^2}{\beta} \mathbf{w} + \mathbf{n} \right)^T \quad (181)$$

Case 6 of 6: If $\|\mathbf{w}\| = 0$ and $\|\mathbf{m}\| = 0$ and $\beta = 0$. Then

$$\mathbf{G} = -\mathbf{v} \mathbf{v}^+ \mathbf{A}^+ - \mathbf{A}^+ \mathbf{n} \mathbf{n}^+ + \mathbf{v}^+ \mathbf{A}^+ \mathbf{n} \mathbf{v} \mathbf{n}^+ \quad (182)$$

$$= -\frac{1}{\|\mathbf{v}\|^2} \mathbf{v} \mathbf{v}^T \mathbf{A}^+ - \frac{1}{\|\mathbf{n}\|^2} \mathbf{A}^+ \mathbf{n} \mathbf{n}^T + \frac{\mathbf{v}^T \mathbf{A}^+ \mathbf{n}}{\|\mathbf{v}\|^2 \|\mathbf{n}\|^2} \mathbf{v} \mathbf{n}^T \quad (183)$$

3.3 Implication on Inverses

$$\text{If } (\mathbf{A} + \mathbf{B})^{-1} = \mathbf{A}^{-1} + \mathbf{B}^{-1} \quad \text{then} \quad \mathbf{A} \mathbf{B}^{-1} \mathbf{A} = \mathbf{B} \mathbf{A}^{-1} \mathbf{B} \quad (184)$$

See [25].

3.3.1 A PosDef identity

Assume \mathbf{P}, \mathbf{R} to be positive definite and invertible, then

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} = \mathbf{P} \mathbf{B}^T (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})^{-1} \quad (185)$$

See [30].

3.4 Approximations

The following identity is known as the *Neuman series* of a matrix, which holds when $|\lambda_i| < 1$ for all eigenvalues λ_i

$$(\mathbf{I} - \mathbf{A})^{-1} = \sum_{n=0}^{\infty} \mathbf{A}^n \quad (186)$$

which is equivalent to

$$(\mathbf{I} + \mathbf{A})^{-1} = \sum_{n=0}^{\infty} (-1)^n \mathbf{A}^n \quad (187)$$

When $|\lambda_i| < 1$ for all eigenvalues λ_i , it holds that $\mathbf{A} \rightarrow 0$ for $n \rightarrow \infty$, and the following approximations holds

$$(\mathbf{I} - \mathbf{A})^{-1} \cong \mathbf{I} + \mathbf{A} + \mathbf{A}^2 \quad (188)$$

$$(\mathbf{I} + \mathbf{A})^{-1} \cong \mathbf{I} - \mathbf{A} + \mathbf{A}^2 \quad (189)$$

The following approximation is from [22] and holds when \mathbf{A} large and symmetric

$$\mathbf{A} - \mathbf{A}(\mathbf{I} + \mathbf{A})^{-1}\mathbf{A} \cong \mathbf{I} - \mathbf{A}^{-1} \quad (190)$$

If σ^2 is small compared to \mathbf{Q} and \mathbf{M} then

$$(\mathbf{Q} + \sigma^2\mathbf{M})^{-1} \cong \mathbf{Q}^{-1} - \sigma^2\mathbf{Q}^{-1}\mathbf{M}\mathbf{Q}^{-1} \quad (191)$$

Proof:

$$(\mathbf{Q} + \sigma^2\mathbf{M})^{-1} = \quad (192)$$

$$(\mathbf{Q}\mathbf{Q}^{-1}\mathbf{Q} + \sigma^2\mathbf{M}\mathbf{Q}^{-1}\mathbf{Q})^{-1} = \quad (193)$$

$$((\mathbf{I} + \sigma^2\mathbf{M}\mathbf{Q}^{-1})\mathbf{Q})^{-1} = \quad (194)$$

$$\mathbf{Q}^{-1}(\mathbf{I} + \sigma^2\mathbf{M}\mathbf{Q}^{-1})^{-1} \quad (195)$$

This can be rewritten using the Taylor expansion:

$$\mathbf{Q}^{-1}(\mathbf{I} + \sigma^2\mathbf{M}\mathbf{Q}^{-1})^{-1} = \quad (196)$$

$$\mathbf{Q}^{-1}(\mathbf{I} - \sigma^2\mathbf{M}\mathbf{Q}^{-1} + (\sigma^2\mathbf{M}\mathbf{Q}^{-1})^2 - \dots) \cong \mathbf{Q}^{-1} - \sigma^2\mathbf{Q}^{-1}\mathbf{M}\mathbf{Q}^{-1} \quad (197)$$

3.5 Generalized Inverse

3.5.1 Definition

A generalized inverse matrix of the matrix \mathbf{A} is any matrix \mathbf{A}^- such that (see [26])

$$\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A} \quad (198)$$

The matrix \mathbf{A}^- is not unique.

3.6 Pseudo Inverse

3.6.1 Definition

The pseudo inverse (or Moore-Penrose inverse) of a matrix \mathbf{A} is the matrix \mathbf{A}^+ that fulfils

$$\text{I} \quad \mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$$

$$\text{II} \quad \mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$$

$$\text{III} \quad \mathbf{A}\mathbf{A}^+ \text{ symmetric}$$

$$\text{IV} \quad \mathbf{A}^+\mathbf{A} \text{ symmetric}$$

The matrix \mathbf{A}^+ is unique and does always exist. Note that in case of complex matrices, the symmetric condition is substituted by a condition of being Hermitian.

3.6.2 Properties

Assume \mathbf{A}^+ to be the pseudo-inverse of \mathbf{A} , then (See [3] for some of them)

$$(\mathbf{A}^+)^+ = \mathbf{A} \quad (199)$$

$$(\mathbf{A}^T)^+ = (\mathbf{A}^+)^T \quad (200)$$

$$(\mathbf{A}^H)^+ = (\mathbf{A}^+)^H \quad (201)$$

$$(\mathbf{A}^*)^+ = (\mathbf{A}^+)^* \quad (202)$$

$$(\mathbf{A}^+ \mathbf{A}) \mathbf{A}^H = \mathbf{A}^H \quad (203)$$

$$(\mathbf{A}^+ \mathbf{A}) \mathbf{A}^T \neq \mathbf{A}^T \quad (204)$$

$$(c\mathbf{A})^+ = (1/c)\mathbf{A}^+ \quad (205)$$

$$\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^+ \mathbf{A}^T \quad (206)$$

$$\mathbf{A}^+ = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^+ \quad (207)$$

$$(\mathbf{A}^T \mathbf{A})^+ = \mathbf{A}^+ (\mathbf{A}^T)^+ \quad (208)$$

$$(\mathbf{A} \mathbf{A}^T)^+ = (\mathbf{A}^T)^+ \mathbf{A}^+ \quad (209)$$

$$\mathbf{A}^+ = (\mathbf{A}^H \mathbf{A})^+ \mathbf{A}^H \quad (210)$$

$$\mathbf{A}^+ = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^+ \quad (211)$$

$$(\mathbf{A}^H \mathbf{A})^+ = \mathbf{A}^+ (\mathbf{A}^H)^+ \quad (212)$$

$$(\mathbf{A} \mathbf{A}^H)^+ = (\mathbf{A}^H)^+ \mathbf{A}^+ \quad (213)$$

$$(\mathbf{A} \mathbf{B})^+ = (\mathbf{A}^+ \mathbf{A} \mathbf{B})^+ (\mathbf{A} \mathbf{B} \mathbf{B}^+)^+ \quad (214)$$

$$f(\mathbf{A}^H \mathbf{A}) - f(0)\mathbf{I} = \mathbf{A}^+ [f(\mathbf{A} \mathbf{A}^H) - f(0)\mathbf{I}] \mathbf{A} \quad (215)$$

$$f(\mathbf{A} \mathbf{A}^H) - f(0)\mathbf{I} = \mathbf{A} [f(\mathbf{A}^H \mathbf{A}) - f(0)\mathbf{I}] \mathbf{A}^+ \quad (216)$$

where $\mathbf{A} \in \mathbb{C}^{n \times m}$.

Assume \mathbf{A} to have full rank, then

$$(\mathbf{A} \mathbf{A}^+)(\mathbf{A} \mathbf{A}^+) = \mathbf{A} \mathbf{A}^+ \quad (217)$$

$$(\mathbf{A}^+ \mathbf{A})(\mathbf{A}^+ \mathbf{A}) = \mathbf{A}^+ \mathbf{A} \quad (218)$$

$$\text{Tr}(\mathbf{A} \mathbf{A}^+) = \text{rank}(\mathbf{A} \mathbf{A}^+) \quad (\text{See [26]}) \quad (219)$$

$$\text{Tr}(\mathbf{A}^+ \mathbf{A}) = \text{rank}(\mathbf{A}^+ \mathbf{A}) \quad (\text{See [26]}) \quad (220)$$

For two matrices it hold that

$$(\mathbf{A} \mathbf{B})^+ = (\mathbf{A}^+ \mathbf{A} \mathbf{B})^+ (\mathbf{A} \mathbf{B} \mathbf{B}^+)^+ \quad (221)$$

$$(\mathbf{A} \otimes \mathbf{B})^+ = \mathbf{A}^+ \otimes \mathbf{B}^+ \quad (222)$$

3.6.3 Construction

Assume that \mathbf{A} has full rank, then

$\mathbf{A} \ n \times n$	Square	$\text{rank}(\mathbf{A}) = n$	\Rightarrow	$\mathbf{A}^+ = \mathbf{A}^{-1}$
$\mathbf{A} \ n \times m$	Broad	$\text{rank}(\mathbf{A}) = n$	\Rightarrow	$\mathbf{A}^+ = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1}$
$\mathbf{A} \ n \times m$	Tall	$\text{rank}(\mathbf{A}) = m$	\Rightarrow	$\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$

The so-called "broad version" is also known as *right inverse* and the "tall version" as the *left inverse*.

Assume \mathbf{A} does not have full rank, i.e. \mathbf{A} is $n \times m$ and $\text{rank}(\mathbf{A}) = r < \min(n, m)$. The pseudo inverse \mathbf{A}^+ can be constructed from the singular value decomposition $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$, by

$$\mathbf{A}^+ = \mathbf{V}_r \mathbf{D}_r^{-1} \mathbf{U}_r^T \quad (223)$$

where \mathbf{U}_r , \mathbf{D}_r , and \mathbf{V}_r are the matrices with the degenerated rows and columns deleted. A different way is this: There do always exist two matrices \mathbf{C} $n \times r$ and \mathbf{D} $r \times m$ of rank r , such that $\mathbf{A} = \mathbf{C}\mathbf{D}$. Using these matrices it holds that

$$\mathbf{A}^+ = \mathbf{D}^T (\mathbf{D}\mathbf{D}^T)^{-1} (\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \quad (224)$$

See [3].

4 Complex Matrices

The complex scalar product $r = pq$ can be written as

$$\begin{bmatrix} \Re r \\ \Im r \end{bmatrix} = \begin{bmatrix} \Re p & -\Im p \\ \Im p & \Re p \end{bmatrix} \begin{bmatrix} \Re q \\ \Im q \end{bmatrix} \quad (225)$$

4.1 Complex Derivatives

In order to differentiate an expression $f(z)$ with respect to a complex z , the Cauchy-Riemann equations have to be satisfied ([7]):

$$\frac{df(z)}{dz} = \frac{\partial \Re(f(z))}{\partial \Re z} + i \frac{\partial \Im(f(z))}{\partial \Re z} \quad (226)$$

and

$$\frac{df(z)}{dz} = -i \frac{\partial \Re(f(z))}{\partial \Im z} + \frac{\partial \Im(f(z))}{\partial \Im z} \quad (227)$$

or in a more compact form:

$$\frac{\partial f(z)}{\partial \Im z} = i \frac{\partial f(z)}{\partial \Re z}. \quad (228)$$

A complex function that satisfies the Cauchy-Riemann equations for points in a region R is said to be *analytic* in this region R . In general, expressions involving complex conjugate or conjugate transpose do not satisfy the Cauchy-Riemann equations. In order to avoid this problem, a more generalized definition of complex derivative is used ([24], [6]):

- Generalized Complex Derivative:

$$\frac{df(z)}{dz} = \frac{1}{2} \left(\frac{\partial f(z)}{\partial \Re z} - i \frac{\partial f(z)}{\partial \Im z} \right). \quad (229)$$

- Conjugate Complex Derivative

$$\frac{df(z)}{dz^*} = \frac{1}{2} \left(\frac{\partial f(z)}{\partial \Re z} + i \frac{\partial f(z)}{\partial \Im z} \right). \quad (230)$$

The Generalized Complex Derivative equals the normal derivative, when f is an analytic function. For a non-analytic function such as $f(z) = z^*$, the derivative equals zero. The Conjugate Complex Derivative equals zero, when f is an analytic function. The Conjugate Complex Derivative has e.g. been used by [21] when deriving a complex gradient.

Notice:

$$\frac{df(z)}{dz} \neq \frac{\partial f(z)}{\partial \Re z} + i \frac{\partial f(z)}{\partial \Im z}. \quad (231)$$

- Complex Gradient Vector: If f is a real function of a complex vector \mathbf{z} , then the complex gradient vector is given by ([14, p. 798])

$$\begin{aligned} \nabla f(\mathbf{z}) &= 2 \frac{df(\mathbf{z})}{d\mathbf{z}^*} \\ &= \frac{\partial f(\mathbf{z})}{\partial \Re \mathbf{z}} + i \frac{\partial f(\mathbf{z})}{\partial \Im \mathbf{z}}. \end{aligned} \quad (232)$$

- **Complex Gradient Matrix:** If f is a real function of a complex matrix \mathbf{Z} , then the complex gradient matrix is given by ([2])

$$\begin{aligned}\nabla f(\mathbf{Z}) &= 2 \frac{df(\mathbf{Z})}{d\mathbf{Z}^*} \\ &= \frac{\partial f(\mathbf{Z})}{\partial \Re \mathbf{Z}} + i \frac{\partial f(\mathbf{Z})}{\partial \Im \mathbf{Z}}.\end{aligned}\quad (233)$$

These expressions can be used for gradient descent algorithms.

4.1.1 The Chain Rule for complex numbers

The chain rule is a little more complicated when the function of a complex $u = f(x)$ is non-analytic. For a non-analytic function, the following chain rule can be applied ([7])

$$\begin{aligned}\frac{\partial g(u)}{\partial x} &= \frac{\partial g}{\partial u} \frac{\partial u}{\partial x} + \frac{\partial g}{\partial u^*} \frac{\partial u^*}{\partial x} \\ &= \frac{\partial g}{\partial u} \frac{\partial u}{\partial x} + \left(\frac{\partial g}{\partial u} \right)^* \frac{\partial u^*}{\partial x}\end{aligned}\quad (234)$$

Notice, if the function is analytic, the second term reduces to zero, and the function is reduced to the normal well-known chain rule. For the matrix derivative of a scalar function $g(\mathbf{U})$, the chain rule can be written the following way:

$$\frac{\partial g(\mathbf{U})}{\partial \mathbf{X}} = \frac{\text{Tr}((\frac{\partial g(\mathbf{U})}{\partial \mathbf{U}})^T \partial \mathbf{U})}{\partial \mathbf{X}} + \frac{\text{Tr}((\frac{\partial g(\mathbf{U})}{\partial \mathbf{U}^*})^T \partial \mathbf{U}^*)}{\partial \mathbf{X}}. \quad (235)$$

4.1.2 Complex Derivatives of Traces

If the derivatives involve complex numbers, the conjugate transpose is often involved. The most useful way to show complex derivative is to show the derivative with respect to the real and the imaginary part separately. An easy example is:

$$\frac{\partial \text{Tr}(\mathbf{X}^*)}{\partial \Re \mathbf{X}} = \frac{\partial \text{Tr}(\mathbf{X}^H)}{\partial \Re \mathbf{X}} = \mathbf{I} \quad (236)$$

$$i \frac{\partial \text{Tr}(\mathbf{X}^*)}{\partial \Im \mathbf{X}} = i \frac{\partial \text{Tr}(\mathbf{X}^H)}{\partial \Im \mathbf{X}} = \mathbf{I} \quad (237)$$

Since the two results have the same sign, the conjugate complex derivative (230) should be used.

$$\frac{\partial \text{Tr}(\mathbf{X})}{\partial \Re \mathbf{X}} = \frac{\partial \text{Tr}(\mathbf{X}^T)}{\partial \Re \mathbf{X}} = \mathbf{I} \quad (238)$$

$$i \frac{\partial \text{Tr}(\mathbf{X})}{\partial \Im \mathbf{X}} = i \frac{\partial \text{Tr}(\mathbf{X}^T)}{\partial \Im \mathbf{X}} = -\mathbf{I} \quad (239)$$

Here, the two results have different signs, and the generalized complex derivative (229) should be used. Hereby, it can be seen that (100) holds even if \mathbf{X} is a complex number.

$$\frac{\partial \text{Tr}(\mathbf{A}\mathbf{X}^H)}{\partial \Re \mathbf{X}} = \mathbf{A} \quad (240)$$

$$i \frac{\partial \text{Tr}(\mathbf{A}\mathbf{X}^H)}{\partial \Im \mathbf{X}} = \mathbf{A} \quad (241)$$

$$\frac{\partial \text{Tr}(\mathbf{A}\mathbf{X}^*)}{\partial \Re \mathbf{X}} = \mathbf{A}^T \quad (242)$$

$$i \frac{\partial \text{Tr}(\mathbf{A}\mathbf{X}^*)}{\partial \Im \mathbf{X}} = \mathbf{A}^T \quad (243)$$

$$\frac{\partial \text{Tr}(\mathbf{X}\mathbf{X}^H)}{\partial \Re \mathbf{X}} = \frac{\partial \text{Tr}(\mathbf{X}^H \mathbf{X})}{\partial \Re \mathbf{X}} = 2\Re \mathbf{X} \quad (244)$$

$$i \frac{\partial \text{Tr}(\mathbf{X}\mathbf{X}^H)}{\partial \Im \mathbf{X}} = i \frac{\partial \text{Tr}(\mathbf{X}^H \mathbf{X})}{\partial \Im \mathbf{X}} = i2\Im \mathbf{X} \quad (245)$$

By inserting (244) and (245) in (229) and (230), it can be seen that

$$\frac{\partial \text{Tr}(\mathbf{X}\mathbf{X}^H)}{\partial \mathbf{X}} = \mathbf{X}^* \quad (246)$$

$$\frac{\partial \text{Tr}(\mathbf{X}\mathbf{X}^H)}{\partial \mathbf{X}^*} = \mathbf{X} \quad (247)$$

Since the function $\text{Tr}(\mathbf{X}\mathbf{X}^H)$ is a real function of the complex matrix \mathbf{X} , the complex gradient matrix (233) is given by

$$\nabla \text{Tr}(\mathbf{X}\mathbf{X}^H) = 2 \frac{\partial \text{Tr}(\mathbf{X}\mathbf{X}^H)}{\partial \mathbf{X}^*} = 2\mathbf{X} \quad (248)$$

4.1.3 Complex Derivative Involving Determinants

Here, a calculation example is provided. The objective is to find the derivative of $\det(\mathbf{X}^H \mathbf{A} \mathbf{X})$ with respect to $\mathbf{X} \in \mathbb{C}^{m \times n}$. The derivative is found with respect to the real part and the imaginary part of \mathbf{X} , by use of (42) and (37), $\det(\mathbf{X}^H \mathbf{A} \mathbf{X})$ can be calculated as (see App. B.1.4 for details)

$$\begin{aligned} \frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \mathbf{X}} &= \frac{1}{2} \left(\frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Re \mathbf{X}} - i \frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Im \mathbf{X}} \right) \\ &= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) ((\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \mathbf{A})^T \end{aligned} \quad (249)$$

and the complex conjugate derivative yields

$$\begin{aligned} \frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \mathbf{X}^*} &= \frac{1}{2} \left(\frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Re \mathbf{X}} + i \frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Im \mathbf{X}} \right) \\ &= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) \mathbf{A} \mathbf{X} (\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \end{aligned} \quad (250)$$

4.2 Higher order and non-linear derivatives

$$\frac{\partial (\mathbf{A}\mathbf{x})^H (\mathbf{A}\mathbf{x})}{\partial \mathbf{x} (\mathbf{B}\mathbf{x})^H (\mathbf{B}\mathbf{x})} = \frac{\partial \mathbf{x}^H \mathbf{A}^H \mathbf{A} \mathbf{x}}{\partial \mathbf{x} \mathbf{x}^H \mathbf{B}^H \mathbf{B} \mathbf{x}} \quad (251)$$

$$= 2 \frac{\mathbf{A}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{B} \mathbf{B} \mathbf{x}} - 2 \frac{\mathbf{x}^H \mathbf{A}^H \mathbf{A} \mathbf{x} \mathbf{B}^H \mathbf{B} \mathbf{x}}{(\mathbf{x}^H \mathbf{B}^H \mathbf{B} \mathbf{x})^2} \quad (252)$$

4.3 Inverse of complex sum

Given real matrices \mathbf{A}, \mathbf{B} find the inverse of the complex sum $\mathbf{A} + i\mathbf{B}$. Form the auxiliary matrices

$$\mathbf{E} = \mathbf{A} + t\mathbf{B} \quad (253)$$

$$\mathbf{F} = \mathbf{B} - t\mathbf{A}, \quad (254)$$

and find a value of t such that \mathbf{E}^{-1} exists. Then

$$(\mathbf{A} + i\mathbf{B})^{-1} = (1 - it)(\mathbf{E} + i\mathbf{F})^{-1} \quad (255)$$

$$= (1 - it)((\mathbf{E} + \mathbf{F}\mathbf{E}^{-1}\mathbf{F})^{-1} - i(\mathbf{E} + \mathbf{F}\mathbf{E}^{-1}\mathbf{F})^{-1}\mathbf{F}\mathbf{E}^{-1}) \quad (256)$$

$$= (1 - it)(\mathbf{E} + \mathbf{F}\mathbf{E}^{-1}\mathbf{F})^{-1}(\mathbf{I} - i\mathbf{F}\mathbf{E}^{-1}) \quad (257)$$

$$= (\mathbf{E} + \mathbf{F}\mathbf{E}^{-1}\mathbf{F})^{-1}((\mathbf{I} - t\mathbf{F}\mathbf{E}^{-1}) - i(t\mathbf{I} + \mathbf{F}\mathbf{E}^{-1})) \quad (258)$$

$$= (\mathbf{E} + \mathbf{F}\mathbf{E}^{-1}\mathbf{F})^{-1}(\mathbf{I} - t\mathbf{F}\mathbf{E}^{-1}) - i(\mathbf{E} + \mathbf{F}\mathbf{E}^{-1}\mathbf{F})^{-1}(t\mathbf{I} + \mathbf{F}\mathbf{E}^{-1}) \quad (259)$$

5 Solutions and Decompositions

5.1 Solutions to linear equations

5.1.1 Simple Linear Regression

Assume we have data (x_n, y_n) for $n = 1, \dots, N$ and are seeking the parameters $a, b \in \mathbb{R}$ such that $y_i \cong ax_i + b$. With a least squares error function, the optimal values for a, b can be expressed using the notation

$$\mathbf{x} = (x_1, \dots, x_N)^T \quad \mathbf{y} = (y_1, \dots, y_N)^T \quad \mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^{N \times 1}$$

and

$$\begin{aligned} R_{xx} &= \mathbf{x}^T \mathbf{x} & R_{x1} &= \mathbf{x}^T \mathbf{1} & R_{11} &= \mathbf{1}^T \mathbf{1} \\ R_{yx} &= \mathbf{y}^T \mathbf{x} & R_{y1} &= \mathbf{y}^T \mathbf{1} \end{aligned}$$

as

$$\begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} R_{xx} & R_{x1} \\ R_{x1} & R_{11} \end{bmatrix}^{-1} \begin{bmatrix} R_{x,y} \\ R_{y1} \end{bmatrix} \quad (260)$$

5.1.2 Existence in Linear Systems

Assume \mathbf{A} is $n \times m$ and consider the linear system

$$\mathbf{Ax} = \mathbf{b} \quad (261)$$

Construct the augmented matrix $\mathbf{B} = [\mathbf{A} \ \mathbf{b}]$ then

<i>Condition</i>	<i>Solution</i>
$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{B}) = m$	Unique solution \mathbf{x}
$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{B}) < m$	Many solutions \mathbf{x}
$\text{rank}(\mathbf{A}) < \text{rank}(\mathbf{B})$	No solutions \mathbf{x}

5.1.3 Standard Square

Assume \mathbf{A} is square and invertible, then

$$\mathbf{Ax} = \mathbf{b} \quad \Rightarrow \quad \mathbf{x} = \mathbf{A}^{-1}\mathbf{b} \quad (262)$$

5.1.4 Degenerated Square

Assume \mathbf{A} is $n \times n$ but of rank $r < n$. In that case, the system $\mathbf{Ax} = \mathbf{b}$ is solved by

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

where \mathbf{A}^+ is the pseudo-inverse of the rank-deficient matrix, constructed as described in section 3.6.3.

5.1.5 Cramer's rule

The equation

$$\mathbf{Ax} = \mathbf{b}, \quad (263)$$

where \mathbf{A} is square has exactly one solution \mathbf{x} if the i th element in \mathbf{x} can be found as

$$x_i = \frac{\det \mathbf{B}}{\det \mathbf{A}}, \quad (264)$$

where \mathbf{B} equals \mathbf{A} , but the i th column in \mathbf{A} has been substituted by \mathbf{b} .

5.1.6 Over-determined Rectangular

Assume \mathbf{A} to be $n \times m$, $n > m$ (tall) and $\text{rank}(\mathbf{A}) = m$, then

$$\mathbf{Ax} = \mathbf{b} \quad \Rightarrow \quad \mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} = \mathbf{A}^+ \mathbf{b} \quad (265)$$

that is *if* there exists a solution \mathbf{x} at all! If there is no solution the following can be useful:

$$\mathbf{Ax} = \mathbf{b} \quad \Rightarrow \quad \mathbf{x}_{min} = \mathbf{A}^+ \mathbf{b} \quad (266)$$

Now \mathbf{x}_{min} is the vector \mathbf{x} which minimizes $\|\mathbf{Ax} - \mathbf{b}\|^2$, i.e. the vector which is "least wrong". The matrix \mathbf{A}^+ is the pseudo-inverse of \mathbf{A} . See [3].

5.1.7 Under-determined Rectangular

Assume \mathbf{A} is $n \times m$ and $n < m$ ("broad") and $\text{rank}(\mathbf{A}) = n$.

$$\mathbf{Ax} = \mathbf{b} \quad \Rightarrow \quad \mathbf{x}_{min} = \mathbf{A}^T (\mathbf{AA}^T)^{-1} \mathbf{b} \quad (267)$$

The equation have many solutions \mathbf{x} . But \mathbf{x}_{min} is the solution which minimizes $\|\mathbf{Ax} - \mathbf{b}\|^2$ and also the solution with the smallest norm $\|\mathbf{x}\|^2$. The same holds for a matrix version: Assume \mathbf{A} is $n \times m$, \mathbf{X} is $m \times n$ and \mathbf{B} is $n \times n$, then

$$\mathbf{AX} = \mathbf{B} \quad \Rightarrow \quad \mathbf{X}_{min} = \mathbf{A}^+ \mathbf{B} \quad (268)$$

The equation have many solutions \mathbf{X} . But \mathbf{X}_{min} is the solution which minimizes $\|\mathbf{AX} - \mathbf{B}\|^2$ and also the solution with the smallest norm $\|\mathbf{X}\|^2$. See [3].

Similar but different: Assume \mathbf{A} is square $n \times n$ and the matrices $\mathbf{B}_0, \mathbf{B}_1$ are $n \times N$, where $N > n$, then if \mathbf{B}_0 has maximal rank

$$\mathbf{AB}_0 = \mathbf{B}_1 \quad \Rightarrow \quad \mathbf{A}_{min} = \mathbf{B}_1 \mathbf{B}_0^T (\mathbf{B}_0 \mathbf{B}_0^T)^{-1} \quad (269)$$

where \mathbf{A}_{min} denotes the matrix which is optimal in a least square sense. An interpretation is that \mathbf{A} is the linear approximation which maps the columns vectors of \mathbf{B}_0 into the columns vectors of \mathbf{B}_1 .

5.1.8 Linear form and zeros

$$\mathbf{Ax} = \mathbf{0}, \quad \forall \mathbf{x} \quad \Rightarrow \quad \mathbf{A} = \mathbf{0} \quad (270)$$

5.1.9 Square form and zeros

If \mathbf{A} is symmetric, then

$$\mathbf{x}^T \mathbf{Ax} = 0, \quad \forall \mathbf{x} \quad \Rightarrow \quad \mathbf{A} = \mathbf{0} \quad (271)$$

5.1.10 The Lyapunov Equation

$$\mathbf{A}\mathbf{X} + \mathbf{X}\mathbf{B} = \mathbf{C} \quad (272)$$

$$\text{vec}(\mathbf{X}) = (\mathbf{I} \otimes \mathbf{A} + \mathbf{B}^T \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{C}) \quad (273)$$

Sec 10.2.1 and 10.2.2 for details on the Kronecker product and the vec operator.

5.1.11 Encapsulating Sum

$$\sum_n \mathbf{A}_n \mathbf{X} \mathbf{B}_n = \mathbf{C} \quad (274)$$

$$\text{vec}(\mathbf{X}) = \left(\sum_n \mathbf{B}_n^T \otimes \mathbf{A}_n \right)^{-1} \text{vec}(\mathbf{C}) \quad (275)$$

See Sec 10.2.1 and 10.2.2 for details on the Kronecker product and the vec operator.

5.2 Eigenvalues and Eigenvectors

5.2.1 Definition

The eigenvectors \mathbf{v}_i and eigenvalues λ_i are the ones satisfying

$$\mathbf{A}\mathbf{v}_i = \lambda_i \mathbf{v}_i \quad (276)$$

5.2.2 Decompositions

For matrices \mathbf{A} with as many distinct eigenvalues as dimensions, the following holds, where the columns of \mathbf{V} are the eigenvectors and $(\mathbf{D})_{ij} = \delta_{ij} \lambda_i$,

$$\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{D} \quad (277)$$

For defective matrices \mathbf{A} , which is matrices which has fewer distinct eigenvalues than dimensions, the following decomposition called *Jordan canonical form*, holds

$$\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{J} \quad (278)$$

where \mathbf{J} is a block diagonal matrix with the blocks $\mathbf{J}_i = \lambda_i \mathbf{I} + \mathbf{N}$. The matrices \mathbf{J}_i have dimensionality as the number of identical eigenvalues equal to λ_i , and \mathbf{N} is square matrix of same size with 1 on the super diagonal and zero elsewhere.

It also holds that for all matrices \mathbf{A} there exists matrices \mathbf{V} and \mathbf{R} such that

$$\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{R} \quad (279)$$

where \mathbf{R} is upper triangular with the eigenvalues λ_i on its diagonal.

5.2.3 General Properties

Assume that $\mathbf{A} \in \mathbb{R}^{n \times m}$ and $\mathbf{B} \in \mathbb{R}^{m \times n}$,

$$\text{eig}(\mathbf{A}\mathbf{B}) = \text{eig}(\mathbf{B}\mathbf{A}) \quad (280)$$

$$\text{rank}(\mathbf{A}) = r \Rightarrow \text{At most } r \text{ non-zero } \lambda_i \quad (281)$$

5.2.4 Symmetric

Assume \mathbf{A} is symmetric, then

$$\mathbf{V}\mathbf{V}^T = \mathbf{I} \quad (\text{i.e. } \mathbf{V} \text{ is orthogonal}) \quad (282)$$

$$\lambda_i \in \mathbb{R} \quad (\text{i.e. } \lambda_i \text{ is real}) \quad (283)$$

$$\text{Tr}(\mathbf{A}^p) = \sum_i \lambda_i^p \quad (284)$$

$$\text{eig}(\mathbf{I} + c\mathbf{A}) = 1 + c\lambda_i \quad (285)$$

$$\text{eig}(\mathbf{A} - c\mathbf{I}) = \lambda_i - c \quad (286)$$

$$\text{eig}(\mathbf{A}^{-1}) = \lambda_i^{-1} \quad (287)$$

For a symmetric, positive matrix \mathbf{A} ,

$$\text{eig}(\mathbf{A}^T \mathbf{A}) = \text{eig}(\mathbf{A}\mathbf{A}^T) = \text{eig}(\mathbf{A}) \circ \text{eig}(\mathbf{A}) \quad (288)$$

5.2.5 Characteristic polynomial

The characteristic polynomial for the matrix \mathbf{A} is

$$0 = \det(\mathbf{A} - \lambda\mathbf{I}) \quad (289)$$

$$= \lambda^n - g_1\lambda^{n-1} + g_2\lambda^{n-2} - \dots + (-1)^n g_n \quad (290)$$

Note that the coefficients g_j for $j = 1, \dots, n$ are the n invariants under rotation of \mathbf{A} . Thus, g_j is the sum of the determinants of all the sub-matrices of \mathbf{A} taken j rows and columns at a time. That is, g_1 is the trace of \mathbf{A} , and g_2 is the sum of the determinants of the $n(n-1)/2$ sub-matrices that can be formed from \mathbf{A} by deleting all but two rows and columns, and so on – see [17].

5.3 Singular Value Decomposition

Any $n \times m$ matrix \mathbf{A} can be written as

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T, \quad (291)$$

where

$$\begin{aligned} \mathbf{U} &= \text{eigenvectors of } \mathbf{A}\mathbf{A}^T & n \times n \\ \mathbf{D} &= \sqrt{\text{diag}(\text{eig}(\mathbf{A}\mathbf{A}^T))} & n \times m \\ \mathbf{V} &= \text{eigenvectors of } \mathbf{A}^T \mathbf{A} & m \times m \end{aligned} \quad (292)$$

5.3.1 Symmetric Square decomposed into squares

Assume \mathbf{A} to be $n \times n$ and symmetric. Then

$$[\mathbf{A}] = [\mathbf{V}] [\mathbf{D}] [\mathbf{V}^T], \quad (293)$$

where \mathbf{D} is diagonal with the eigenvalues of \mathbf{A} , and \mathbf{V} is orthogonal and the eigenvectors of \mathbf{A} .

5.3.2 Square decomposed into squares

Assume $\mathbf{A} \in \mathbb{R}^{n \times n}$. Then

$$[\mathbf{A}] = [\mathbf{V}] [\mathbf{D}] [\mathbf{U}^T], \quad (294)$$

where \mathbf{D} is diagonal with the square root of the eigenvalues of $\mathbf{A}\mathbf{A}^T$, \mathbf{V} is the eigenvectors of $\mathbf{A}\mathbf{A}^T$ and \mathbf{U}^T is the eigenvectors of $\mathbf{A}^T \mathbf{A}$.

5.3.3 Square decomposed into rectangular

Assume $\mathbf{V}_* \mathbf{D}_* \mathbf{U}_*^T = \mathbf{0}$ then we can expand the SVD of \mathbf{A} into

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{V} & \mathbf{V}_* \end{bmatrix} \left[\begin{array}{c|c} \mathbf{D} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{D}_* \end{array} \right] \begin{bmatrix} \mathbf{U}^T \\ \mathbf{U}_*^T \end{bmatrix}, \quad (295)$$

where the SVD of \mathbf{A} is $\mathbf{A} = \mathbf{V} \mathbf{D} \mathbf{U}^T$.

5.3.4 Rectangular decomposition I

Assume \mathbf{A} is $n \times m$, \mathbf{V} is $n \times n$, \mathbf{D} is $n \times n$, \mathbf{U}^T is $n \times m$

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{V} \end{bmatrix} \begin{bmatrix} \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{U}^T \end{bmatrix}, \quad (296)$$

where \mathbf{D} is diagonal with the square root of the eigenvalues of $\mathbf{A} \mathbf{A}^T$, \mathbf{V} is the eigenvectors of $\mathbf{A} \mathbf{A}^T$ and \mathbf{U}^T is the eigenvectors of $\mathbf{A}^T \mathbf{A}$.

5.3.5 Rectangular decomposition II

Assume \mathbf{A} is $n \times m$, \mathbf{V} is $n \times m$, \mathbf{D} is $m \times m$, \mathbf{U}^T is $m \times m$

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{V} \end{bmatrix} \begin{bmatrix} \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{U}^T \end{bmatrix} \quad (297)$$

5.3.6 Rectangular decomposition III

Assume \mathbf{A} is $n \times m$, \mathbf{V} is $n \times n$, \mathbf{D} is $n \times m$, \mathbf{U}^T is $m \times m$

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{V} \end{bmatrix} \begin{bmatrix} \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{U}^T \end{bmatrix}, \quad (298)$$

where \mathbf{D} is diagonal with the square root of the eigenvalues of $\mathbf{A} \mathbf{A}^T$, \mathbf{V} is the eigenvectors of $\mathbf{A} \mathbf{A}^T$ and \mathbf{U}^T is the eigenvectors of $\mathbf{A}^T \mathbf{A}$.

5.4 Triangular Decomposition

5.5 LU decomposition

Assume \mathbf{A} is a square matrix with non-zero leading principal minors, then

$$\mathbf{A} = \mathbf{L} \mathbf{U} \quad (299)$$

where \mathbf{L} is a unique unit lower triangular matrix and \mathbf{U} is a unique upper triangular matrix.

5.5.1 Cholesky-decomposition

Assume \mathbf{A} is a symmetric positive definite square matrix, then

$$\mathbf{A} = \mathbf{U}^T \mathbf{U} = \mathbf{L} \mathbf{L}^T, \quad (300)$$

where \mathbf{U} is a unique upper triangular matrix and \mathbf{L} is a lower triangular matrix.

5.6 LDM decomposition

Assume \mathbf{A} is a square matrix with non-zero leading principal minors¹, then

$$\mathbf{A} = \mathbf{LDM}^T \quad (301)$$

where \mathbf{L} , \mathbf{M} are unique unit lower triangular matrices and \mathbf{D} is a unique diagonal matrix.

5.7 LDL decompositions

The LDL decomposition are special cases of the LDM decomposition. Assume \mathbf{A} is a non-singular symmetric definite square matrix, then

$$\mathbf{A} = \mathbf{LDL}^T = \mathbf{L}^T \mathbf{D} \mathbf{L} \quad (302)$$

where \mathbf{L} is a unit lower triangular matrix and \mathbf{D} is a diagonal matrix. If \mathbf{A} is also positive definite, then \mathbf{D} has strictly positive diagonal entries.

¹If the matrix that corresponds to a principal minor is a quadratic upper-left part of the larger matrix (i.e., it consists of matrix elements in rows and columns from 1 to k), then the principal minor is called a leading principal minor. For an n times n square matrix, there are n leading principal minors. [31]

6 Statistics and Probability

6.1 Definition of Moments

Assume $\mathbf{x} \in \mathbb{R}^{n \times 1}$ is a random variable

6.1.1 Mean

The vector of means, \mathbf{m} , is defined by

$$(\mathbf{m})_i = \langle x_i \rangle \quad (303)$$

6.1.2 Covariance

The matrix of covariance \mathbf{M} is defined by

$$(\mathbf{M})_{ij} = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \quad (304)$$

or alternatively as

$$\mathbf{M} = \langle (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T \rangle \quad (305)$$

6.1.3 Third moments

The matrix of third centralized moments – in some contexts referred to as coskewness – is defined using the notation

$$m_{ijk}^{(3)} = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle)(x_k - \langle x_k \rangle) \rangle \quad (306)$$

as

$$\mathbf{M}_3 = \left[m_{::1}^{(3)} m_{::2}^{(3)} \dots m_{::n}^{(3)} \right] \quad (307)$$

where ‘:’ denotes all elements within the given index. \mathbf{M}_3 can alternatively be expressed as

$$\mathbf{M}_3 = \langle (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T \otimes (\mathbf{x} - \mathbf{m})^T \rangle \quad (308)$$

6.1.4 Fourth moments

The matrix of fourth centralized moments – in some contexts referred to as cokurtosis – is defined using the notation

$$m_{ijkl}^{(4)} = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle)(x_k - \langle x_k \rangle)(x_l - \langle x_l \rangle) \rangle \quad (309)$$

as

$$\mathbf{M}_4 = \left[m_{::11}^{(4)} m_{::21}^{(4)} \dots m_{::n1}^{(4)} | m_{::12}^{(4)} m_{::22}^{(4)} \dots m_{::n2}^{(4)} | \dots | m_{::1n}^{(4)} m_{::2n}^{(4)} \dots m_{::nn}^{(4)} \right] \quad (310)$$

or alternatively as

$$\mathbf{M}_4 = \langle (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T \otimes (\mathbf{x} - \mathbf{m})^T \otimes (\mathbf{x} - \mathbf{m})^T \rangle \quad (311)$$

6.2 Expectation of Linear Combinations

6.2.1 Linear Forms

Assume \mathbf{X} and \mathbf{x} to be a matrix and a vector of random variables. Then (see See [26])

$$E[\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C}] = \mathbf{A}E[\mathbf{X}]\mathbf{B} + \mathbf{C} \quad (312)$$

$$\text{Var}[\mathbf{A}\mathbf{x}] = \mathbf{A}\text{Var}[\mathbf{x}]\mathbf{A}^T \quad (313)$$

$$\text{Cov}[\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y}] = \mathbf{A}\text{Cov}[\mathbf{x}, \mathbf{y}]\mathbf{B}^T \quad (314)$$

Assume \mathbf{x} to be a stochastic vector with mean \mathbf{m} , then (see [7])

$$E[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbf{m} + \mathbf{b} \quad (315)$$

$$E[\mathbf{A}\mathbf{x}] = \mathbf{A}\mathbf{m} \quad (316)$$

$$E[\mathbf{x} + \mathbf{b}] = \mathbf{m} + \mathbf{b} \quad (317)$$

6.2.2 Quadratic Forms

Assume \mathbf{A} is symmetric, $\mathbf{c} = E[\mathbf{x}]$ and $\mathbf{\Sigma} = \text{Var}[\mathbf{x}]$. Assume also that all coordinates x_i are independent, have the same central moments $\mu_1, \mu_2, \mu_3, \mu_4$ and denote $\mathbf{a} = \text{diag}(\mathbf{A})$. Then (See [26])

$$E[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \text{Tr}(\mathbf{A}\mathbf{\Sigma}) + \mathbf{c}^T \mathbf{A} \mathbf{c} \quad (318)$$

$$\text{Var}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = 2\mu_2^2 \text{Tr}(\mathbf{A}^2) + 4\mu_2 \mathbf{c}^T \mathbf{A}^2 \mathbf{c} + 4\mu_3 \mathbf{c}^T \mathbf{A} \mathbf{a} + (\mu_4 - 3\mu_2^2) \mathbf{a}^T \mathbf{a} \quad (319)$$

Also, assume \mathbf{x} to be a stochastic vector with mean \mathbf{m} , and covariance \mathbf{M} . Then (see [7])

$$E[(\mathbf{A}\mathbf{x} + \mathbf{a})(\mathbf{B}\mathbf{x} + \mathbf{b})^T] = \mathbf{A}\mathbf{M}\mathbf{B}^T + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{B}\mathbf{m} + \mathbf{b})^T \quad (320)$$

$$E[\mathbf{x}\mathbf{x}^T] = \mathbf{M} + \mathbf{m}\mathbf{m}^T \quad (321)$$

$$E[\mathbf{x}\mathbf{a}^T \mathbf{x}] = (\mathbf{M} + \mathbf{m}\mathbf{m}^T)\mathbf{a} \quad (322)$$

$$E[\mathbf{x}^T \mathbf{a}\mathbf{x}^T] = \mathbf{a}^T(\mathbf{M} + \mathbf{m}\mathbf{m}^T) \quad (323)$$

$$E[(\mathbf{A}\mathbf{x})(\mathbf{A}\mathbf{x})^T] = \mathbf{A}(\mathbf{M} + \mathbf{m}\mathbf{m}^T)\mathbf{A}^T \quad (324)$$

$$E[(\mathbf{x} + \mathbf{a})(\mathbf{x} + \mathbf{a})^T] = \mathbf{M} + (\mathbf{m} + \mathbf{a})(\mathbf{m} + \mathbf{a})^T \quad (325)$$

$$E[(\mathbf{A}\mathbf{x} + \mathbf{a})^T(\mathbf{B}\mathbf{x} + \mathbf{b})] = \text{Tr}(\mathbf{A}\mathbf{M}\mathbf{B}^T) + (\mathbf{A}\mathbf{m} + \mathbf{a})^T(\mathbf{B}\mathbf{m} + \mathbf{b}) \quad (326)$$

$$E[\mathbf{x}^T \mathbf{x}] = \text{Tr}(\mathbf{M}) + \mathbf{m}^T \mathbf{m} \quad (327)$$

$$E[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \text{Tr}(\mathbf{A}\mathbf{M}) + \mathbf{m}^T \mathbf{A} \mathbf{m} \quad (328)$$

$$E[(\mathbf{A}\mathbf{x})^T(\mathbf{A}\mathbf{x})] = \text{Tr}(\mathbf{A}\mathbf{M}\mathbf{A}^T) + (\mathbf{A}\mathbf{m})^T(\mathbf{A}\mathbf{m}) \quad (329)$$

$$E[(\mathbf{x} + \mathbf{a})^T(\mathbf{x} + \mathbf{a})] = \text{Tr}(\mathbf{M}) + (\mathbf{m} + \mathbf{a})^T(\mathbf{m} + \mathbf{a}) \quad (330)$$

See [7].

6.2.3 Cubic Forms

Assume \mathbf{x} to be a stochastic vector with independent coordinates, mean \mathbf{m} , covariance \mathbf{M} and central moments $\mathbf{v}_3 = E[(\mathbf{x} - \mathbf{m})^3]$. Then (see [7])

$$\begin{aligned}
E[(\mathbf{Ax} + \mathbf{a})(\mathbf{Bx} + \mathbf{b})^T(\mathbf{Cx} + \mathbf{c})] &= \mathbf{A} \text{diag}(\mathbf{B}^T \mathbf{C}) \mathbf{v}_3 \\
&\quad + \text{Tr}(\mathbf{BMC}^T)(\mathbf{Am} + \mathbf{a}) \\
&\quad + \mathbf{AMC}^T(\mathbf{Bm} + \mathbf{b}) \\
&\quad + (\mathbf{AMB}^T + (\mathbf{Am} + \mathbf{a})(\mathbf{Bm} + \mathbf{b})^T)(\mathbf{Cm} + \mathbf{c}) \\
E[\mathbf{xx}^T \mathbf{x}] &= \mathbf{v}_3 + 2\mathbf{Mm} + (\text{Tr}(\mathbf{M}) + \mathbf{m}^T \mathbf{m}) \mathbf{m} \\
E[(\mathbf{Ax} + \mathbf{a})(\mathbf{Ax} + \mathbf{a})^T(\mathbf{Ax} + \mathbf{a})] &= \mathbf{A} \text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{v}_3 \\
&\quad + [2\mathbf{AMA}^T + (\mathbf{Ax} + \mathbf{a})(\mathbf{Ax} + \mathbf{a})^T](\mathbf{Am} + \mathbf{a}) \\
&\quad + \text{Tr}(\mathbf{AMA}^T)(\mathbf{Am} + \mathbf{a}) \\
E[(\mathbf{Ax} + \mathbf{a})\mathbf{b}^T(\mathbf{Cx} + \mathbf{c})(\mathbf{Dx} + \mathbf{d})^T] &= (\mathbf{Ax} + \mathbf{a})\mathbf{b}^T(\mathbf{CMD}^T + (\mathbf{Cm} + \mathbf{c})(\mathbf{Dm} + \mathbf{d})^T) \\
&\quad + (\mathbf{AMC}^T + (\mathbf{Am} + \mathbf{a})(\mathbf{Cm} + \mathbf{c})^T)\mathbf{b}(\mathbf{Dm} + \mathbf{d})^T \\
&\quad + \mathbf{b}^T(\mathbf{Cm} + \mathbf{c})(\mathbf{AMD}^T - (\mathbf{Am} + \mathbf{a})(\mathbf{Dm} + \mathbf{d})^T)
\end{aligned}$$

6.3 Weighted Scalar Variable

Assume $\mathbf{x} \in \mathbb{R}^{n \times 1}$ is a random variable, $\mathbf{w} \in \mathbb{R}^{n \times 1}$ is a vector of constants and y is the linear combination $y = \mathbf{w}^T \mathbf{x}$. Assume further that $\mathbf{m}, \mathbf{M}_2, \mathbf{M}_3, \mathbf{M}_4$ denotes the mean, covariance, and central third and fourth moment matrix of the variable \mathbf{x} . Then it holds that

$$\langle y \rangle = \mathbf{w}^T \mathbf{m} \quad (331)$$

$$\langle (y - \langle y \rangle)^2 \rangle = \mathbf{w}^T \mathbf{M}_2 \mathbf{w} \quad (332)$$

$$\langle (y - \langle y \rangle)^3 \rangle = \mathbf{w}^T \mathbf{M}_3 \mathbf{w} \otimes \mathbf{w} \quad (333)$$

$$\langle (y - \langle y \rangle)^4 \rangle = \mathbf{w}^T \mathbf{M}_4 \mathbf{w} \otimes \mathbf{w} \otimes \mathbf{w} \quad (334)$$

7 Multivariate Distributions

7.1 Cauchy

The density function for a Cauchy distributed vector $\mathbf{t} \in \mathbb{R}^{P \times 1}$, is given by

$$p(\mathbf{t}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \pi^{-P/2} \frac{\Gamma(\frac{1+P}{2})}{\Gamma(1/2)} \frac{\det(\boldsymbol{\Sigma})^{-1/2}}{[1 + (\mathbf{t} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{t} - \boldsymbol{\mu})]^{(1+P)/2}} \quad (335)$$

where $\boldsymbol{\mu}$ is the location, $\boldsymbol{\Sigma}$ is positive definite, and Γ denotes the gamma function. The Cauchy distribution is a special case of the Student-t distribution.

7.2 Dirichlet

The Dirichlet distribution is a kind of “inverse” distribution compared to the multinomial distribution on the bounded continuous variate $\mathbf{x} = [x_1, \dots, x_P]$ [16, p. 44]

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_p \alpha_p)}{\prod_p \Gamma(\alpha_p)} \prod_p x_p^{\alpha_p-1}$$

7.3 Normal

The normal distribution is also known as a Gaussian distribution. See sec. 8.

7.4 Normal-Inverse Gamma

7.5 Gaussian

See sec. 8.

7.6 Multinomial

If the vector \mathbf{n} contains counts, i.e. $(\mathbf{n})_i \in 0, 1, 2, \dots$, then the discrete multinomial distribution for \mathbf{n} is given by

$$P(\mathbf{n}|\mathbf{a}, n) = \frac{n!}{n_1! \dots n_d!} \prod_i^d a_i^{n_i}, \quad \sum_i^d n_i = n \quad (336)$$

where a_i are probabilities, i.e. $0 \leq a_i \leq 1$ and $\sum_i a_i = 1$.

7.7 Student's t

The density of a Student-t distributed vector $\mathbf{t} \in \mathbb{R}^{P \times 1}$, is given by

$$p(\mathbf{t}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) = (\pi\nu)^{-P/2} \frac{\Gamma(\frac{\nu+P}{2})}{\Gamma(\nu/2)} \frac{\det(\boldsymbol{\Sigma})^{-1/2}}{[1 + \nu^{-1}(\mathbf{t} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{t} - \boldsymbol{\mu})]^{(\nu+P)/2}} \quad (337)$$

where $\boldsymbol{\mu}$ is the location, the scale matrix $\boldsymbol{\Sigma}$ is symmetric, positive definite, ν is the degrees of freedom, and Γ denotes the gamma function. For $\nu = 1$, the Student-t distribution becomes the Cauchy distribution (see sec 7.1).

7.7.1 Mean

$$E(\mathbf{t}) = \boldsymbol{\mu}, \quad \nu > 1 \quad (338)$$

7.7.2 Variance

$$\text{cov}(\mathbf{t}) = \frac{\nu}{\nu - 2} \boldsymbol{\Sigma}, \quad \nu > 2 \quad (339)$$

7.7.3 Mode

The notion *mode* meaning the position of the most probable value

$$\text{mode}(\mathbf{t}) = \boldsymbol{\mu} \quad (340)$$

7.7.4 Full Matrix Version

If instead of a vector $\mathbf{t} \in \mathbb{R}^{P \times 1}$ one has a matrix $\mathbf{T} \in \mathbb{R}^{P \times N}$, then the Student-t distribution for \mathbf{T} is

$$\begin{aligned} p(\mathbf{T}|\mathbf{M}, \boldsymbol{\Omega}, \boldsymbol{\Sigma}, \nu) &= \pi^{-NP/2} \prod_{p=1}^P \frac{\Gamma[(\nu + P - p + 1)/2]}{\Gamma[(\nu - p + 1)/2]} \times \\ &\quad \nu \det(\boldsymbol{\Omega})^{-\nu/2} \det(\boldsymbol{\Sigma})^{-N/2} \times \\ &\quad \det[\boldsymbol{\Omega}^{-1} + (\mathbf{T} - \mathbf{M})\boldsymbol{\Sigma}^{-1}(\mathbf{T} - \mathbf{M})^T]^{-(\nu+P)/2} \end{aligned} \quad (341)$$

where \mathbf{M} is the location, $\boldsymbol{\Omega}$ is the rescaling matrix, $\boldsymbol{\Sigma}$ is positive definite, ν is the degrees of freedom, and Γ denotes the gamma function.

7.8 Wishart

The central Wishart distribution for $\mathbf{M} \in \mathbb{R}^{P \times P}$, \mathbf{M} is positive definite, where m can be regarded as a degree of freedom parameter [16, equation 3.8.1] [8, section 2.5],[11]

$$\begin{aligned} p(\mathbf{M}|\boldsymbol{\Sigma}, m) &= \frac{1}{2^{mP/2} \pi^{P(P-1)/4} \prod_p^P \Gamma[\frac{1}{2}(m+1-p)]} \times \\ &\quad \det(\boldsymbol{\Sigma})^{-m/2} \det(\mathbf{M})^{(m-P-1)/2} \times \\ &\quad \exp\left[-\frac{1}{2}\text{Tr}(\boldsymbol{\Sigma}^{-1}\mathbf{M})\right] \end{aligned} \quad (342)$$

7.8.1 Mean

$$E(\mathbf{M}) = \mathbf{m}\boldsymbol{\Sigma} \quad (343)$$

7.9 Wishart, Inverse

The (normal) Inverse Wishart distribution for $\mathbf{M} \in \mathbb{R}^{P \times P}$, \mathbf{M} is positive definite, where m can be regarded as a degree of freedom parameter [11]

$$\begin{aligned}
 p(\mathbf{M}|\boldsymbol{\Sigma}, m) &= \frac{1}{2^{mP/2} \pi^{P(P-1)/4} \prod_p^P \Gamma[\frac{1}{2}(m+1-p)]} \times \\
 &\quad \det(\boldsymbol{\Sigma})^{m/2} \det(\mathbf{M})^{-(m-P+1)/2} \times \\
 &\quad \exp \left[-\frac{1}{2} \text{Tr}(\boldsymbol{\Sigma} \mathbf{M}^{-1}) \right]
 \end{aligned} \tag{344}$$

7.9.1 Mean

$$E(\mathbf{M}) = \boldsymbol{\Sigma} \frac{1}{m - P + 1} \tag{345}$$

8 Gaussians

8.1 Basics

8.1.1 Density and normalization

The density of $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$ is

$$p(\mathbf{x}) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp \left[-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \Sigma^{-1}(\mathbf{x} - \mathbf{m}) \right] \quad (346)$$

Note that if \mathbf{x} is d -dimensional, then $\det(2\pi\Sigma) = (2\pi)^d \det(\Sigma)$.
Integration and normalization

$$\begin{aligned} \int \exp \left[-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \Sigma^{-1}(\mathbf{x} - \mathbf{m}) \right] d\mathbf{x} &= \sqrt{\det(2\pi\Sigma)} \\ \int \exp \left[-\frac{1}{2}\mathbf{x}^T \Sigma^{-1}\mathbf{x} + \mathbf{m}^T \Sigma^{-1}\mathbf{x} \right] d\mathbf{x} &= \sqrt{\det(2\pi\Sigma)} \exp \left[\frac{1}{2}\mathbf{m}^T \Sigma^{-1}\mathbf{m} \right] \\ \int \exp \left[-\frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + \mathbf{c}^T \mathbf{x} \right] d\mathbf{x} &= \sqrt{\det(2\pi\mathbf{A}^{-1})} \exp \left[\frac{1}{2}\mathbf{c}^T \mathbf{A}^{-1}\mathbf{c} \right] \end{aligned}$$

If $\mathbf{X} = [\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_n]$ and $\mathbf{C} = [\mathbf{c}_1 \mathbf{c}_2 \dots \mathbf{c}_n]$, then

$$\int \exp \left[-\frac{1}{2}\text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X}) + \text{Tr}(\mathbf{C}^T \mathbf{X}) \right] d\mathbf{X} = \sqrt{\det(2\pi\mathbf{A}^{-1})}^n \exp \left[\frac{1}{2}\text{Tr}(\mathbf{C}^T \mathbf{A}^{-1} \mathbf{C}) \right]$$

The derivatives of the density are

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{x}} = -p(\mathbf{x}) \Sigma^{-1}(\mathbf{x} - \mathbf{m}) \quad (347)$$

$$\frac{\partial^2 p}{\partial \mathbf{x} \partial \mathbf{x}^T} = p(\mathbf{x}) \left(\Sigma^{-1}(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T \Sigma^{-1} - \Sigma^{-1} \right) \quad (348)$$

8.1.2 Marginal Distribution

Assume $\mathbf{x} \sim \mathcal{N}_{\mathbf{x}}(\mu, \Sigma)$ where

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_a & \Sigma_c \\ \Sigma_c^T & \Sigma_b \end{bmatrix} \quad (349)$$

then

$$p(\mathbf{x}_a) = \mathcal{N}_{\mathbf{x}_a}(\mu_a, \Sigma_a) \quad (350)$$

$$p(\mathbf{x}_b) = \mathcal{N}_{\mathbf{x}_b}(\mu_b, \Sigma_b) \quad (351)$$

8.1.3 Conditional Distribution

Assume $\mathbf{x} \sim \mathcal{N}_{\mathbf{x}}(\mu, \Sigma)$ where

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_a & \Sigma_c \\ \Sigma_c^T & \Sigma_b \end{bmatrix} \quad (352)$$

then

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}_{\mathbf{x}_a}(\hat{\boldsymbol{\mu}}_a, \hat{\boldsymbol{\Sigma}}_a) \quad \begin{cases} \hat{\boldsymbol{\mu}}_a &= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_c \boldsymbol{\Sigma}_b^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ \hat{\boldsymbol{\Sigma}}_a &= \boldsymbol{\Sigma}_a - \boldsymbol{\Sigma}_c \boldsymbol{\Sigma}_b^{-1} \boldsymbol{\Sigma}_c^T \end{cases} \quad (353)$$

$$p(\mathbf{x}_b|\mathbf{x}_a) = \mathcal{N}_{\mathbf{x}_b}(\hat{\boldsymbol{\mu}}_b, \hat{\boldsymbol{\Sigma}}_b) \quad \begin{cases} \hat{\boldsymbol{\mu}}_b &= \boldsymbol{\mu}_b + \boldsymbol{\Sigma}_c^T \boldsymbol{\Sigma}_a^{-1}(\mathbf{x}_a - \boldsymbol{\mu}_a) \\ \hat{\boldsymbol{\Sigma}}_b &= \boldsymbol{\Sigma}_b - \boldsymbol{\Sigma}_c^T \boldsymbol{\Sigma}_a^{-1} \boldsymbol{\Sigma}_c \end{cases} \quad (354)$$

Note, that the covariance matrices are the Schur complement of the block matrix, see 9.1.5 for details.

8.1.4 Linear combination

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_x, \boldsymbol{\Sigma}_x)$ and $\mathbf{y} \sim \mathcal{N}(\mathbf{m}_y, \boldsymbol{\Sigma}_y)$ then

$$\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c} \sim \mathcal{N}(\mathbf{A}\mathbf{m}_x + \mathbf{B}\mathbf{m}_y + \mathbf{c}, \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_y\mathbf{B}^T) \quad (355)$$

8.1.5 Rearranging Means

$$\mathcal{N}_{\mathbf{A}\mathbf{x}}[\mathbf{m}, \boldsymbol{\Sigma}] = \frac{\sqrt{\det(2\pi(\mathbf{A}^T\boldsymbol{\Sigma}^{-1}\mathbf{A})^{-1})}}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \mathcal{N}_{\mathbf{x}}[\mathbf{A}^{-1}\mathbf{m}, (\mathbf{A}^T\boldsymbol{\Sigma}^{-1}\mathbf{A})^{-1}] \quad (356)$$

If \mathbf{A} is square and invertible, it simplifies to

$$\mathcal{N}_{\mathbf{A}\mathbf{x}}[\mathbf{m}, \boldsymbol{\Sigma}] = \frac{1}{|\det(\mathbf{A})|} \mathcal{N}_{\mathbf{x}}[\mathbf{A}^{-1}\mathbf{m}, (\mathbf{A}^T\boldsymbol{\Sigma}^{-1}\mathbf{A})^{-1}] \quad (357)$$

8.1.6 Rearranging into squared form

If \mathbf{A} is symmetric, then

$$\begin{aligned} -\frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} + \mathbf{b}^T\mathbf{x} &= -\frac{1}{2}(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b})^T\mathbf{A}(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b}) + \frac{1}{2}\mathbf{b}^T\mathbf{A}^{-1}\mathbf{b} \\ -\frac{1}{2}\text{Tr}(\mathbf{X}^T\mathbf{A}\mathbf{X}) + \text{Tr}(\mathbf{B}^T\mathbf{X}) &= -\frac{1}{2}\text{Tr}[(\mathbf{X} - \mathbf{A}^{-1}\mathbf{B})^T\mathbf{A}(\mathbf{X} - \mathbf{A}^{-1}\mathbf{B})] + \frac{1}{2}\text{Tr}(\mathbf{B}^T\mathbf{A}^{-1}\mathbf{B}) \end{aligned}$$

8.1.7 Sum of two squared forms

In vector formulation (assuming $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2$ are symmetric)

$$-\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^T\boldsymbol{\Sigma}_1^{-1}(\mathbf{x} - \mathbf{m}_1) \quad (358)$$

$$-\frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^T\boldsymbol{\Sigma}_2^{-1}(\mathbf{x} - \mathbf{m}_2) \quad (359)$$

$$= -\frac{1}{2}(\mathbf{x} - \mathbf{m}_c)^T\boldsymbol{\Sigma}_c^{-1}(\mathbf{x} - \mathbf{m}_c) + C \quad (360)$$

$$\boldsymbol{\Sigma}_c^{-1} = \boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1} \quad (361)$$

$$\mathbf{m}_c = (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}(\boldsymbol{\Sigma}_1^{-1}\mathbf{m}_1 + \boldsymbol{\Sigma}_2^{-1}\mathbf{m}_2) \quad (362)$$

$$C = \frac{1}{2}(\mathbf{m}_1^T\boldsymbol{\Sigma}_1^{-1} + \mathbf{m}_2^T\boldsymbol{\Sigma}_2^{-1})(\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}(\boldsymbol{\Sigma}_1^{-1}\mathbf{m}_1 + \boldsymbol{\Sigma}_2^{-1}\mathbf{m}_2) \quad (363)$$

$$-\frac{1}{2}\left(\mathbf{m}_1^T\boldsymbol{\Sigma}_1^{-1}\mathbf{m}_1 + \mathbf{m}_2^T\boldsymbol{\Sigma}_2^{-1}\mathbf{m}_2\right) \quad (364)$$

In a trace formulation (assuming Σ_1, Σ_2 are symmetric)

$$-\frac{1}{2}\text{Tr}((\mathbf{X} - \mathbf{M}_1)^T \Sigma_1^{-1} (\mathbf{X} - \mathbf{M}_1)) \quad (365)$$

$$-\frac{1}{2}\text{Tr}((\mathbf{X} - \mathbf{M}_2)^T \Sigma_2^{-1} (\mathbf{X} - \mathbf{M}_2)) \quad (366)$$

$$= -\frac{1}{2}\text{Tr}[(\mathbf{X} - \mathbf{M}_c)^T \Sigma_c^{-1} (\mathbf{X} - \mathbf{M}_c)] + C \quad (367)$$

$$\Sigma_c^{-1} = \Sigma_1^{-1} + \Sigma_2^{-1} \quad (368)$$

$$\mathbf{M}_c = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} (\Sigma_1^{-1} \mathbf{M}_1 + \Sigma_2^{-1} \mathbf{M}_2) \quad (369)$$

$$\begin{aligned} C &= \frac{1}{2}\text{Tr}[(\Sigma_1^{-1} \mathbf{M}_1 + \Sigma_2^{-1} \mathbf{M}_2)^T (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} (\Sigma_1^{-1} \mathbf{M}_1 + \Sigma_2^{-1} \mathbf{M}_2)] \\ &\quad - \frac{1}{2}\text{Tr}(\mathbf{M}_1^T \Sigma_1^{-1} \mathbf{M}_1 + \mathbf{M}_2^T \Sigma_2^{-1} \mathbf{M}_2) \end{aligned} \quad (370)$$

8.1.8 Product of gaussian densities

Let $\mathcal{N}_{\mathbf{x}}(\mathbf{m}, \Sigma)$ denote a density of \mathbf{x} , then

$$\mathcal{N}_{\mathbf{x}}(\mathbf{m}_1, \Sigma_1) \cdot \mathcal{N}_{\mathbf{x}}(\mathbf{m}_2, \Sigma_2) = c_c \mathcal{N}_{\mathbf{x}}(\mathbf{m}_c, \Sigma_c) \quad (371)$$

$$\begin{aligned} c_c &= \mathcal{N}_{\mathbf{m}_1}(\mathbf{m}_2, (\Sigma_1 + \Sigma_2)) \\ &= \frac{1}{\sqrt{\det(2\pi(\Sigma_1 + \Sigma_2))}} \exp\left[-\frac{1}{2}(\mathbf{m}_1 - \mathbf{m}_2)^T (\Sigma_1 + \Sigma_2)^{-1} (\mathbf{m}_1 - \mathbf{m}_2)\right] \\ \mathbf{m}_c &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} (\Sigma_1^{-1} \mathbf{m}_1 + \Sigma_2^{-1} \mathbf{m}_2) \\ \Sigma_c &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \end{aligned}$$

but note that the product is not normalized as a density of \mathbf{x} .

8.2 Moments

8.2.1 Mean and covariance of linear forms

First and second moments. Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$E(\mathbf{x}) = \mathbf{m} \quad (372)$$

$$\text{Cov}(\mathbf{x}, \mathbf{x}) = \text{Var}(\mathbf{x}) = \Sigma = E(\mathbf{x}\mathbf{x}^T) - E(\mathbf{x})E(\mathbf{x}^T) = E(\mathbf{x}\mathbf{x}^T) - \mathbf{m}\mathbf{m}^T \quad (373)$$

As for any other distribution is holds for gaussians that

$$E[\mathbf{A}\mathbf{x}] = \mathbf{A}E[\mathbf{x}] \quad (374)$$

$$\text{Var}[\mathbf{A}\mathbf{x}] = \mathbf{A}\text{Var}[\mathbf{x}]\mathbf{A}^T \quad (375)$$

$$\text{Cov}[\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y}] = \mathbf{A}\text{Cov}[\mathbf{x}, \mathbf{y}]\mathbf{B}^T \quad (376)$$

8.2.2 Mean and variance of square forms

Mean and variance of square forms: Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$E(\mathbf{x}\mathbf{x}^T) = \Sigma + \mathbf{m}\mathbf{m}^T \quad (377)$$

$$E[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \text{Tr}(\mathbf{A}\Sigma) + \mathbf{m}^T \mathbf{A} \mathbf{m} \quad (378)$$

$$\begin{aligned} \text{Var}(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \text{Tr}[\mathbf{A}\Sigma(\mathbf{A} + \mathbf{A}^T)\Sigma] + \dots \\ &\quad + \mathbf{m}^T (\mathbf{A} + \mathbf{A}^T) \Sigma (\mathbf{A} + \mathbf{A}^T) \mathbf{m} \end{aligned} \quad (379)$$

$$E[(\mathbf{x} - \mathbf{m}')^T \mathbf{A} (\mathbf{x} - \mathbf{m}')] = (\mathbf{m} - \mathbf{m}')^T \mathbf{A} (\mathbf{m} - \mathbf{m}') + \text{Tr}(\mathbf{A}\Sigma) \quad (380)$$

If $\Sigma = \sigma^2 \mathbf{I}$ and \mathbf{A} is symmetric, then

$$\text{Var}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = 2\sigma^4 \text{Tr}(\mathbf{A}^2) + 4\sigma^2 \mathbf{m}^T \mathbf{A}^2 \mathbf{m} \quad (381)$$

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ and \mathbf{A} and \mathbf{B} to be symmetric, then

$$\text{Cov}(\mathbf{x}^T \mathbf{A} \mathbf{x}, \mathbf{x}^T \mathbf{B} \mathbf{x}) = 2\sigma^4 \text{Tr}(\mathbf{A}\mathbf{B}) \quad (382)$$

8.2.3 Cubic forms

Assume \mathbf{x} to be a stochastic vector with independent coordinates, mean \mathbf{m} and covariance \mathbf{M}

$$\begin{aligned} E[\mathbf{x}\mathbf{b}^T \mathbf{x}\mathbf{x}^T] &= \mathbf{m}\mathbf{b}^T (\mathbf{M} + \mathbf{m}\mathbf{m}^T) + (\mathbf{M} + \mathbf{m}\mathbf{m}^T) \mathbf{b}\mathbf{m}^T \\ &\quad + \mathbf{b}^T \mathbf{m} (\mathbf{M} - \mathbf{m}\mathbf{m}^T) \end{aligned} \quad (383)$$

8.2.4 Mean of Quartic Forms

$$\begin{aligned} E[\mathbf{x}\mathbf{x}^T \mathbf{x}\mathbf{x}^T] &= 2(\Sigma + \mathbf{m}\mathbf{m}^T)^2 + \mathbf{m}^T \mathbf{m} (\Sigma - \mathbf{m}\mathbf{m}^T) \\ &\quad + \text{Tr}(\Sigma)(\Sigma + \mathbf{m}\mathbf{m}^T) \\ E[\mathbf{x}\mathbf{x}^T \mathbf{A}\mathbf{x}\mathbf{x}^T] &= (\Sigma + \mathbf{m}\mathbf{m}^T)(\mathbf{A} + \mathbf{A}^T)(\Sigma + \mathbf{m}\mathbf{m}^T) \\ &\quad + \mathbf{m}^T \mathbf{A} \mathbf{m} (\Sigma - \mathbf{m}\mathbf{m}^T) + \text{Tr}[\mathbf{A}\Sigma](\Sigma + \mathbf{m}\mathbf{m}^T) \\ E[\mathbf{x}^T \mathbf{x}\mathbf{x}^T \mathbf{x}] &= 2\text{Tr}(\Sigma^2) + 4\mathbf{m}^T \Sigma \mathbf{m} + (\text{Tr}(\Sigma) + \mathbf{m}^T \mathbf{m})^2 \\ E[\mathbf{x}^T \mathbf{A}\mathbf{x}\mathbf{x}^T \mathbf{B}\mathbf{x}] &= \text{Tr}[\mathbf{A}\Sigma(\mathbf{B} + \mathbf{B}^T)\Sigma] + \mathbf{m}^T (\mathbf{A} + \mathbf{A}^T) \Sigma (\mathbf{B} + \mathbf{B}^T) \mathbf{m} \\ &\quad + (\text{Tr}(\mathbf{A}\Sigma) + \mathbf{m}^T \mathbf{A} \mathbf{m})(\text{Tr}(\mathbf{B}\Sigma) + \mathbf{m}^T \mathbf{B} \mathbf{m}) \end{aligned}$$

$$\begin{aligned} &E[\mathbf{a}^T \mathbf{x}\mathbf{b}^T \mathbf{x}\mathbf{c}^T \mathbf{x}\mathbf{d}^T \mathbf{x}] \\ &= (\mathbf{a}^T (\Sigma + \mathbf{m}\mathbf{m}^T) \mathbf{b})(\mathbf{c}^T (\Sigma + \mathbf{m}\mathbf{m}^T) \mathbf{d}) \\ &\quad + (\mathbf{a}^T (\Sigma + \mathbf{m}\mathbf{m}^T) \mathbf{c})(\mathbf{b}^T (\Sigma + \mathbf{m}\mathbf{m}^T) \mathbf{d}) \\ &\quad + (\mathbf{a}^T (\Sigma + \mathbf{m}\mathbf{m}^T) \mathbf{d})(\mathbf{b}^T (\Sigma + \mathbf{m}\mathbf{m}^T) \mathbf{c}) - 2\mathbf{a}^T \mathbf{m}\mathbf{b}^T \mathbf{m}\mathbf{c}^T \mathbf{m}\mathbf{d}^T \mathbf{m} \\ &E[(\mathbf{A}\mathbf{x} + \mathbf{a})(\mathbf{B}\mathbf{x} + \mathbf{b})^T (\mathbf{C}\mathbf{x} + \mathbf{c})(\mathbf{D}\mathbf{x} + \mathbf{d})^T] \\ &= [\mathbf{A}\Sigma\mathbf{B}^T + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{B}\mathbf{m} + \mathbf{b})^T][\mathbf{C}\Sigma\mathbf{D}^T + (\mathbf{C}\mathbf{m} + \mathbf{c})(\mathbf{D}\mathbf{m} + \mathbf{d})^T] \\ &\quad + [\mathbf{A}\Sigma\mathbf{C}^T + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{C}\mathbf{m} + \mathbf{c})^T][\mathbf{B}\Sigma\mathbf{D}^T + (\mathbf{B}\mathbf{m} + \mathbf{b})(\mathbf{D}\mathbf{m} + \mathbf{d})^T] \\ &\quad + (\mathbf{B}\mathbf{m} + \mathbf{b})^T (\mathbf{C}\mathbf{m} + \mathbf{c})[\mathbf{A}\Sigma\mathbf{D}^T - (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{D}\mathbf{m} + \mathbf{d})^T] \\ &\quad + \text{Tr}(\mathbf{B}\Sigma\mathbf{C}^T)[\mathbf{A}\Sigma\mathbf{D}^T + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{D}\mathbf{m} + \mathbf{d})^T] \end{aligned}$$

$$\begin{aligned}
& E[(\mathbf{A}\mathbf{x} + \mathbf{a})^T (\mathbf{B}\mathbf{x} + \mathbf{b})(\mathbf{C}\mathbf{x} + \mathbf{c})^T (\mathbf{D}\mathbf{x} + \mathbf{d})] \\
&= \text{Tr}[\mathbf{A}\boldsymbol{\Sigma}(\mathbf{C}^T \mathbf{D} + \mathbf{D}^T \mathbf{C})\boldsymbol{\Sigma} \mathbf{B}^T] \\
&\quad + [(\mathbf{A}\mathbf{m} + \mathbf{a})^T \mathbf{B} + (\mathbf{B}\mathbf{m} + \mathbf{b})^T \mathbf{A}] \boldsymbol{\Sigma} [\mathbf{C}^T (\mathbf{D}\mathbf{m} + \mathbf{d}) + \mathbf{D}^T (\mathbf{C}\mathbf{m} + \mathbf{c})] \\
&\quad + [\text{Tr}(\mathbf{A}\boldsymbol{\Sigma} \mathbf{B}^T) + (\mathbf{A}\mathbf{m} + \mathbf{a})^T (\mathbf{B}\mathbf{m} + \mathbf{b})] [\text{Tr}(\mathbf{C}\boldsymbol{\Sigma} \mathbf{D}^T) + (\mathbf{C}\mathbf{m} + \mathbf{c})^T (\mathbf{D}\mathbf{m} + \mathbf{d})]
\end{aligned}$$

See [7].

8.2.5 Moments

$$E[\mathbf{x}] = \sum_k \rho_k \mathbf{m}_k \quad (384)$$

$$\text{Cov}(\mathbf{x}) = \sum_k \sum_{k'} \rho_k \rho_{k'} (\boldsymbol{\Sigma}_k + \mathbf{m}_k \mathbf{m}_k^T - \mathbf{m}_k \mathbf{m}_{k'}^T) \quad (385)$$

8.3 Miscellaneous

8.3.1 Whitening

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma})$ then

$$\mathbf{z} = \boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \mathbf{m}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (386)$$

Conversely having $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ one can generate data $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma})$ by setting

$$\mathbf{x} = \boldsymbol{\Sigma}^{1/2} \mathbf{z} + \mathbf{m} \sim \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma}) \quad (387)$$

Note that $\boldsymbol{\Sigma}^{1/2}$ means the matrix which fulfils $\boldsymbol{\Sigma}^{1/2} \boldsymbol{\Sigma}^{1/2} = \boldsymbol{\Sigma}$, and that it exists and is unique since $\boldsymbol{\Sigma}$ is positive definite.

8.3.2 The Chi-Square connection

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma})$ and \mathbf{x} to be n dimensional, then

$$z = (\mathbf{x} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \mathbf{m}) \sim \chi_n^2 \quad (388)$$

where χ_n^2 denotes the Chi square distribution with n degrees of freedom.

8.3.3 Entropy

Entropy of a D -dimensional gaussian

$$H(\mathbf{x}) = - \int \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma}) \ln \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma}) d\mathbf{x} = \ln \sqrt{\det(2\pi \boldsymbol{\Sigma})} + \frac{D}{2} \quad (389)$$

8.4 Mixture of Gaussians

8.4.1 Density

The variable \mathbf{x} is distributed as a mixture of gaussians if it has the density

$$p(\mathbf{x}) = \sum_{k=1}^K \rho_k \frac{1}{\sqrt{\det(2\pi \boldsymbol{\Sigma}_k)}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mathbf{m}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \mathbf{m}_k) \right] \quad (390)$$

where ρ_k sum to 1 and the $\boldsymbol{\Sigma}_k$ all are positive definite.

8.4.2 Derivatives

Defining $p(\mathbf{s}) = \sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ one get

$$\frac{\partial \ln p(\mathbf{s})}{\partial \rho_j} = \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \frac{\partial}{\partial \rho_j} \ln[\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)] \quad (391)$$

$$= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \frac{1}{\rho_j} \quad (392)$$

$$\frac{\partial \ln p(\mathbf{s})}{\partial \boldsymbol{\mu}_j} = \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \frac{\partial}{\partial \boldsymbol{\mu}_j} \ln[\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)] \quad (393)$$

$$= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} [\boldsymbol{\Sigma}_j^{-1}(\mathbf{s} - \boldsymbol{\mu}_j)] \quad (394)$$

$$\frac{\partial \ln p(\mathbf{s})}{\partial \boldsymbol{\Sigma}_j} = \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \frac{\partial}{\partial \boldsymbol{\Sigma}_j} \ln[\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)] \quad (395)$$

$$= \frac{\rho_j \mathcal{N}_s(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \rho_k \mathcal{N}_s(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \frac{1}{2} [-\boldsymbol{\Sigma}_j^{-1} + \boldsymbol{\Sigma}_j^{-1}(\mathbf{s} - \boldsymbol{\mu}_j)(\mathbf{s} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1}] \quad (396)$$

But ρ_k and $\boldsymbol{\Sigma}_k$ needs to be constrained.

9 Special Matrices

9.1 Block matrices

Let \mathbf{A}_{ij} denote the ij th block of \mathbf{A} .

9.1.1 Multiplication

Assuming the dimensions of the blocks matches we have

$$\left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right] \left[\begin{array}{c|c} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \hline \mathbf{B}_{21} & \mathbf{B}_{22} \end{array} \right] = \left[\begin{array}{c|c} \mathbf{A}_{11}\mathbf{B}_{11} + \mathbf{A}_{12}\mathbf{B}_{21} & \mathbf{A}_{11}\mathbf{B}_{12} + \mathbf{A}_{12}\mathbf{B}_{22} \\ \hline \mathbf{A}_{21}\mathbf{B}_{11} + \mathbf{A}_{22}\mathbf{B}_{21} & \mathbf{A}_{21}\mathbf{B}_{12} + \mathbf{A}_{22}\mathbf{B}_{22} \end{array} \right]$$

9.1.2 The Determinant

The determinant can be expressed as by the use of

$$\mathbf{C}_1 = \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} \quad (397)$$

$$\mathbf{C}_2 = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12} \quad (398)$$

as

$$\det \left(\left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right] \right) = \det(\mathbf{A}_{22}) \cdot \det(\mathbf{C}_1) = \det(\mathbf{A}_{11}) \cdot \det(\mathbf{C}_2)$$

9.1.3 The Inverse

The inverse can be expressed as by the use of

$$\mathbf{C}_1 = \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} \quad (399)$$

$$\mathbf{C}_2 = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12} \quad (400)$$

as

$$\begin{aligned} \left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right]^{-1} &= \left[\begin{array}{c|c} \mathbf{C}_1^{-1} & -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{C}_2^{-1} \\ \hline -\mathbf{C}_2^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{C}_2^{-1} \end{array} \right] \\ &= \left[\begin{array}{c|c} \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{C}_2^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & -\mathbf{C}_1^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \hline -\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{C}_1^{-1} & \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{C}_1^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \end{array} \right] \end{aligned}$$

9.1.4 Block diagonal

For block diagonal matrices we have

$$\left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{A}_{22} \end{array} \right]^{-1} = \left[\begin{array}{c|c} (\mathbf{A}_{11})^{-1} & \mathbf{0} \\ \hline \mathbf{0} & (\mathbf{A}_{22})^{-1} \end{array} \right] \quad (401)$$

$$\det \left(\left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{A}_{22} \end{array} \right] \right) = \det(\mathbf{A}_{11}) \cdot \det(\mathbf{A}_{22}) \quad (402)$$

9.1.5 Schur complement

Regard the matrix

$$\left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right]$$

The Schur complement of block \mathbf{A}_{11} of the matrix above is the matrix (denoted \mathbf{C}_2 in the text above)

$$\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$$

The Schur complement of block \mathbf{A}_{22} of the matrix above is the matrix (denoted \mathbf{C}_1 in the text above)

$$\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$$

Using the Schur complement, one can rewrite the inverse of a block matrix

$$\begin{aligned} & \left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right]^{-1} \\ = & \left[\begin{array}{c|c} \mathbf{I} & \mathbf{0} \\ \hline -\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{I} \end{array} \right] \left[\begin{array}{c|c} (\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{A}_{22}^{-1} \end{array} \right] \left[\begin{array}{c|c} \mathbf{I} & -\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \hline \mathbf{0} & \mathbf{I} \end{array} \right] \end{aligned}$$

The Schur complement is useful when solving linear systems of the form

$$\left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right] \left[\begin{array}{c} \mathbf{x}_1 \\ \mathbf{x}_2 \end{array} \right] = \left[\begin{array}{c} \mathbf{b}_1 \\ \mathbf{b}_2 \end{array} \right]$$

which has the following equation for \mathbf{x}_1

$$(\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})\mathbf{x}_1 = \mathbf{b}_1 - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{b}_2$$

When the appropriate inverses exists, this can be solved for \mathbf{x}_1 which can then be inserted in the equation for \mathbf{x}_2 to solve for \mathbf{x}_2 .

9.2 Discrete Fourier Transform Matrix, The

The DFT matrix is an $N \times N$ symmetric matrix \mathbf{W}_N , where the k, n th element is given by

$$W_N^{kn} = e^{-j\frac{2\pi kn}{N}} \quad (403)$$

Thus the discrete Fourier transform (DFT) can be expressed as

$$X(k) = \sum_{n=0}^{N-1} x(n)W_N^{kn}. \quad (404)$$

Likewise the inverse discrete Fourier transform (IDFT) can be expressed as

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k)W_N^{-kn}. \quad (405)$$

The DFT of the vector $\mathbf{x} = [x(0), x(1), \dots, x(N-1)]^T$ can be written in matrix form as

$$\mathbf{X} = \mathbf{W}_N \mathbf{x}, \quad (406)$$

where $\mathbf{X} = [X(0), X(1), \dots, x(N-1)]^T$. The IDFT is similarly given as

$$\mathbf{x} = \mathbf{W}_N^{-1} \mathbf{X}. \quad (407)$$

Some properties of \mathbf{W}_N exist:

$$\mathbf{W}_N^{-1} = \frac{1}{N} \mathbf{W}_N^* \quad (408)$$

$$\mathbf{W}_N \mathbf{W}_N^* = N \mathbf{I} \quad (409)$$

$$\mathbf{W}_N^* = \mathbf{W}_N^H \quad (410)$$

If $W_N = e^{-\frac{j2\pi}{N}}$, then [23]

$$W_N^{m+N/2} = -W_N^m \quad (411)$$

Notice, the DFT matrix is a Vandermonde Matrix.

The following important relation between the circulant matrix and the discrete Fourier transform (DFT) exists

$$\mathbf{T}_C = \mathbf{W}_N^{-1} (\mathbf{I} \circ (\mathbf{W}_N \mathbf{t})) \mathbf{W}_N, \quad (412)$$

where $\mathbf{t} = [t_0, t_1, \dots, t_{n-1}]^T$ is the first row of \mathbf{T}_C .

9.3 Hermitian Matrices and skew-Hermitian

A matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ is called *Hermitian* if

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

For real valued matrices, Hermitian and symmetric matrices are equivalent.

$$\mathbf{A} \text{ is Hermitian} \Leftrightarrow \mathbf{x}^H \mathbf{A} \mathbf{x} \in \mathbb{R}, \quad \forall \mathbf{x} \in \mathbb{C}^{n \times 1} \quad (413)$$

$$\mathbf{A} \text{ is Hermitian} \Leftrightarrow \text{eig}(\mathbf{A}) \in \mathbb{R} \quad (414)$$

Note that

$$\mathbf{A} = \mathbf{B} + i\mathbf{C}$$

where \mathbf{B}, \mathbf{C} are hermitian, then

$$\mathbf{B} = \frac{\mathbf{A} + \mathbf{A}^H}{2}, \quad \mathbf{C} = \frac{\mathbf{A} - \mathbf{A}^H}{2i}$$

9.3.1 Skew-Hermitian

A matrix \mathbf{A} is called *skew-hermitian* if

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

For real valued matrices, skew-Hermitian and skew-symmetric matrices are equivalent.

$$\mathbf{A} \text{ Hermitian} \Leftrightarrow i\mathbf{A} \text{ is skew-hermitian} \quad (415)$$

$$\mathbf{A} \text{ skew-Hermitian} \Leftrightarrow \mathbf{x}^H \mathbf{A} \mathbf{y} = -\mathbf{x}^H \mathbf{A}^H \mathbf{y}, \quad \forall \mathbf{x}, \mathbf{y} \quad (416)$$

$$\mathbf{A} \text{ skew-Hermitian} \Rightarrow \text{eig}(\mathbf{A}) = i\lambda, \quad \lambda \in \mathbb{R} \quad (417)$$

9.4 Idempotent Matrices

A matrix \mathbf{A} is idempotent if

$$\mathbf{A}\mathbf{A} = \mathbf{A}$$

Idempotent matrices \mathbf{A} and \mathbf{B} , have the following properties

$$\mathbf{A}^n = \mathbf{A}, \quad \text{for } n = 1, 2, 3, \dots \quad (418)$$

$$\mathbf{I} - \mathbf{A} \quad \text{is idempotent} \quad (419)$$

$$\mathbf{A}^H \quad \text{is idempotent} \quad (420)$$

$$\mathbf{I} - \mathbf{A}^H \quad \text{is idempotent} \quad (421)$$

$$\text{If } \mathbf{A}\mathbf{B} = \mathbf{B}\mathbf{A} \Rightarrow \mathbf{A}\mathbf{B} \quad \text{is idempotent} \quad (422)$$

$$\text{rank}(\mathbf{A}) = \text{Tr}(\mathbf{A}) \quad (423)$$

$$\mathbf{A}(\mathbf{I} - \mathbf{A}) = \mathbf{0} \quad (424)$$

$$(\mathbf{I} - \mathbf{A})\mathbf{A} = \mathbf{0} \quad (425)$$

$$\mathbf{A}^+ = \mathbf{A} \quad (426)$$

$$f(s\mathbf{I} + t\mathbf{A}) = (\mathbf{I} - \mathbf{A})f(s) + \mathbf{A}f(s + t) \quad (427)$$

Note that $\mathbf{A} - \mathbf{I}$ is not necessarily idempotent.

9.4.1 Nilpotent

A matrix \mathbf{A} is nilpotent if

$$\mathbf{A}^2 = \mathbf{0}$$

A nilpotent matrix has the following property:

$$f(s\mathbf{I} + t\mathbf{A}) = \mathbf{I}f(s) + t\mathbf{A}f'(s) \quad (428)$$

9.4.2 Unipotent

A matrix \mathbf{A} is unipotent if

$$\mathbf{A}\mathbf{A} = \mathbf{I}$$

A unipotent matrix has the following property:

$$f(s\mathbf{I} + t\mathbf{A}) = [(\mathbf{I} + \mathbf{A})f(s + t) + (\mathbf{I} - \mathbf{A})f(s - t)]/2 \quad (429)$$

9.5 Orthogonal matrices

If a square matrix \mathbf{Q} is orthogonal, if and only if,

$$\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbf{I}$$

and then \mathbf{Q} has the following properties

- Its eigenvalues are placed on the unit circle.
- Its eigenvectors are unitary, i.e. have length one.
- The inverse of an orthogonal matrix is orthogonal too.

Basic properties for the orthogonal matrix \mathbf{Q}

$$\begin{aligned}\mathbf{Q}^{-1} &= \mathbf{Q}^T \\ \mathbf{Q}^{-T} &= \mathbf{Q} \\ \mathbf{Q}\mathbf{Q}^T &= \mathbf{I} \\ \mathbf{Q}^T\mathbf{Q} &= \mathbf{I} \\ \det(\mathbf{Q}) &= \pm 1\end{aligned}$$

9.5.1 Ortho-Sym

A matrix \mathbf{Q}_+ which simultaneously is orthogonal and symmetric is called an ortho-sym matrix [20]. Hereby

$$\mathbf{Q}_+^T \mathbf{Q}_+ = \mathbf{I} \quad (430)$$

$$\mathbf{Q}_+ = \mathbf{Q}_+^T \quad (431)$$

The powers of an ortho-sym matrix are given by the following rule

$$\mathbf{Q}_+^k = \frac{1 + (-1)^k}{2} \mathbf{I} + \frac{1 + (-1)^{k+1}}{2} \mathbf{Q}_+ \quad (432)$$

$$= \frac{1 + \cos(k\pi)}{2} \mathbf{I} + \frac{1 - \cos(k\pi)}{2} \mathbf{Q}_+ \quad (433)$$

9.5.2 Ortho-Skew

A matrix which simultaneously is orthogonal and antisymmetric is called an ortho-skew matrix [20]. Hereby

$$\mathbf{Q}_-^H \mathbf{Q}_- = \mathbf{I} \quad (434)$$

$$\mathbf{Q}_- = -\mathbf{Q}_-^H \quad (435)$$

The powers of an ortho-skew matrix are given by the following rule

$$\mathbf{Q}_-^k = \frac{i^k + (-i)^k}{2} \mathbf{I} - i \frac{i^k - (-i)^k}{2} \mathbf{Q}_- \quad (436)$$

$$= \cos(k\frac{\pi}{2}) \mathbf{I} + \sin(k\frac{\pi}{2}) \mathbf{Q}_- \quad (437)$$

9.5.3 Decomposition

A square matrix \mathbf{A} can always be written as a sum of a symmetric \mathbf{A}_+ and an antisymmetric matrix \mathbf{A}_-

$$\mathbf{A} = \mathbf{A}_+ + \mathbf{A}_- \quad (438)$$

9.6 Positive Definite and Semi-definite Matrices

9.6.1 Definitions

A matrix \mathbf{A} is positive definite if and only if

$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0, \quad \forall \mathbf{x} \neq \mathbf{0} \quad (439)$$

A matrix \mathbf{A} is positive semi-definite if and only if

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0, \quad \forall \mathbf{x} \quad (440)$$

Note that if \mathbf{A} is positive definite, then \mathbf{A} is also positive semi-definite.

9.6.2 Eigenvalues

The following holds with respect to the eigenvalues:

$$\begin{aligned} \mathbf{A} \text{ pos. def.} & \Leftrightarrow \text{eig}\left(\frac{\mathbf{A}+\mathbf{A}^H}{2}\right) > 0 \\ \mathbf{A} \text{ pos. semi-def.} & \Leftrightarrow \text{eig}\left(\frac{\mathbf{A}+\mathbf{A}^H}{2}\right) \geq 0 \end{aligned} \quad (441)$$

9.6.3 Trace

The following holds with respect to the trace:

$$\begin{aligned} \mathbf{A} \text{ pos. def.} & \Rightarrow \text{Tr}(\mathbf{A}) > 0 \\ \mathbf{A} \text{ pos. semi-def.} & \Rightarrow \text{Tr}(\mathbf{A}) \geq 0 \end{aligned} \quad (442)$$

9.6.4 Inverse

If \mathbf{A} is positive definite, then \mathbf{A} is invertible and \mathbf{A}^{-1} is also positive definite.

9.6.5 Diagonal

If \mathbf{A} is positive definite, then $A_{ii} > 0, \forall i$

9.6.6 Decomposition I

The matrix \mathbf{A} is positive semi-definite of rank $r \Leftrightarrow$ there exists a matrix \mathbf{B} of rank r such that $\mathbf{A} = \mathbf{B}\mathbf{B}^T$

The matrix \mathbf{A} is positive definite \Leftrightarrow there exists an invertible matrix \mathbf{B} such that $\mathbf{A} = \mathbf{B}\mathbf{B}^T$

9.6.7 Decomposition II

Assume \mathbf{A} is an $n \times n$ positive semi-definite, then there exists an $n \times r$ matrix \mathbf{B} of rank r such that $\mathbf{B}^T \mathbf{A} \mathbf{B} = \mathbf{I}$.

9.6.8 Equation with zeros

Assume \mathbf{A} is positive semi-definite, then $\mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{0} \Rightarrow \mathbf{A} \mathbf{X} = \mathbf{0}$

9.6.9 Rank of product

Assume \mathbf{A} is positive definite, then $\text{rank}(\mathbf{B}\mathbf{A}\mathbf{B}^T) = \text{rank}(\mathbf{B})$

9.6.10 Positive definite property

If \mathbf{A} is $n \times n$ positive definite and \mathbf{B} is $r \times n$ of rank r , then $\mathbf{B}\mathbf{A}\mathbf{B}^T$ is positive definite.

9.6.11 Outer Product

If \mathbf{X} is $n \times r$, where $n \leq r$ and $\text{rank}(\mathbf{X}) = n$, then $\mathbf{X}\mathbf{X}^T$ is positive definite.

9.6.12 Small perturbations

If \mathbf{A} is positive definite and \mathbf{B} is symmetric, then $\mathbf{A} - t\mathbf{B}$ is positive definite for sufficiently small t .

9.6.13 Hadamard inequality

If \mathbf{A} is a positive definite or semi-definite matrix, then

$$\det(\mathbf{A}) \leq \prod_i A_{ii}$$

See [15, pp.477]

9.6.14 Hadamard product relation

Assume that $\mathbf{P} = \mathbf{A}\mathbf{A}^T$ and $\mathbf{Q} = \mathbf{B}\mathbf{B}^T$ are semi positive definite matrices, it then holds that

$$\mathbf{P} \circ \mathbf{Q} = \mathbf{R}\mathbf{R}^T$$

where the columns of \mathbf{R} are constructed as follows: $\mathbf{r}_{i+(j-1)N_A} = \mathbf{a}_i \circ \mathbf{b}_j$, for $i = 1, 2, \dots, N_A$ and $j = 1, 2, \dots, N_B$. The result is unpublished, but reported by Pavel Sakov and Craig Bishop.

9.7 Singleentry Matrix, The**9.7.1 Definition**

The single-entry matrix $\mathbf{J}^{ij} \in \mathbb{R}^{n \times n}$ is defined as the matrix which is zero everywhere except in the entry (i, j) in which it is 1. In a 4×4 example one might have

$$\mathbf{J}^{23} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (443)$$

The single-entry matrix is very useful when working with derivatives of expressions involving matrices.

9.7.2 Swap and Zeros

Assume \mathbf{A} to be $n \times m$ and \mathbf{J}^{ij} to be $m \times p$

$$\mathbf{A}\mathbf{J}^{ij} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \dots & \mathbf{A}_i & \dots & \mathbf{0} \end{bmatrix} \quad (444)$$

i.e. an $n \times p$ matrix of zeros with the i .th column of \mathbf{A} in place of the j .th column. Assume \mathbf{A} to be $n \times m$ and \mathbf{J}^{ij} to be $p \times n$

$$\mathbf{J}^{ij}\mathbf{A} = \begin{bmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \\ \mathbf{A}_j \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} \quad (445)$$

i.e. an $p \times m$ matrix of zeros with the j .th row of \mathbf{A} in the place of the i .th row.

9.7.3 Rewriting product of elements

$$A_{ki}B_{jl} = (\mathbf{A}\mathbf{e}_i\mathbf{e}_j^T\mathbf{B})_{kl} = (\mathbf{A}\mathbf{J}^{ij}\mathbf{B})_{kl} \quad (446)$$

$$A_{ik}B_{lj} = (\mathbf{A}^T\mathbf{e}_i\mathbf{e}_j^T\mathbf{B}^T)_{kl} = (\mathbf{A}^T\mathbf{J}^{ij}\mathbf{B}^T)_{kl} \quad (447)$$

$$A_{ik}B_{jl} = (\mathbf{A}^T\mathbf{e}_i\mathbf{e}_j^T\mathbf{B})_{kl} = (\mathbf{A}^T\mathbf{J}^{ij}\mathbf{B})_{kl} \quad (448)$$

$$A_{ki}B_{lj} = (\mathbf{A}\mathbf{e}_i\mathbf{e}_j^T\mathbf{B}^T)_{kl} = (\mathbf{A}\mathbf{J}^{ij}\mathbf{B}^T)_{kl} \quad (449)$$

9.7.4 Properties of the Singleentry Matrix

If $i = j$

$$\mathbf{J}^{ij}\mathbf{J}^{ij} = \mathbf{J}^{ij} \quad (\mathbf{J}^{ij})^T(\mathbf{J}^{ij})^T = \mathbf{J}^{ij}$$

$$\mathbf{J}^{ij}(\mathbf{J}^{ij})^T = \mathbf{J}^{ij} \quad (\mathbf{J}^{ij})^T\mathbf{J}^{ij} = \mathbf{J}^{ij}$$

If $i \neq j$

$$\mathbf{J}^{ij}\mathbf{J}^{ij} = \mathbf{0} \quad (\mathbf{J}^{ij})^T(\mathbf{J}^{ij})^T = \mathbf{0}$$

$$\mathbf{J}^{ij}(\mathbf{J}^{ij})^T = \mathbf{J}^{ii} \quad (\mathbf{J}^{ij})^T\mathbf{J}^{ij} = \mathbf{J}^{jj}$$

9.7.5 The Singleentry Matrix in Scalar Expressions

Assume \mathbf{A} is $n \times m$ and \mathbf{J} is $m \times n$, then

$$\text{Tr}(\mathbf{A}\mathbf{J}^{ij}) = \text{Tr}(\mathbf{J}^{ij}\mathbf{A}) = (\mathbf{A}^T)_{ij} \quad (450)$$

Assume \mathbf{A} is $n \times n$, \mathbf{J} is $n \times m$ and \mathbf{B} is $m \times n$, then

$$\text{Tr}(\mathbf{A}\mathbf{J}^{ij}\mathbf{B}) = (\mathbf{A}^T\mathbf{B}^T)_{ij} \quad (451)$$

$$\text{Tr}(\mathbf{A}\mathbf{J}^{ji}\mathbf{B}) = (\mathbf{B}\mathbf{A})_{ij} \quad (452)$$

$$\text{Tr}(\mathbf{A}\mathbf{J}^{ij}\mathbf{J}^{ij}\mathbf{B}) = \text{diag}(\mathbf{A}^T\mathbf{B}^T)_{ij} \quad (453)$$

Assume \mathbf{A} is $n \times n$, \mathbf{J}^{ij} is $n \times m$ \mathbf{B} is $m \times n$, then

$$\mathbf{x}^T\mathbf{A}\mathbf{J}^{ij}\mathbf{B}\mathbf{x} = (\mathbf{A}^T\mathbf{x}\mathbf{x}^T\mathbf{B}^T)_{ij} \quad (454)$$

$$\mathbf{x}^T\mathbf{A}\mathbf{J}^{ij}\mathbf{J}^{ij}\mathbf{B}\mathbf{x} = \text{diag}(\mathbf{A}^T\mathbf{x}\mathbf{x}^T\mathbf{B}^T)_{ij} \quad (455)$$

9.7.6 Structure Matrices

The structure matrix is defined by

$$\frac{\partial \mathbf{A}}{\partial A_{ij}} = \mathbf{S}^{ij} \quad (456)$$

If \mathbf{A} has no special structure then

$$\mathbf{S}^{ij} = \mathbf{J}^{ij} \quad (457)$$

If \mathbf{A} is symmetric then

$$\mathbf{S}^{ij} = \mathbf{J}^{ij} + \mathbf{J}^{ji} - \mathbf{J}^{ij}\mathbf{J}^{ij} \quad (458)$$

9.8 Symmetric, Skew-symmetric/Antisymmetric

9.8.1 Symmetric

The matrix \mathbf{A} is said to be *symmetric* if

$$\mathbf{A} = \mathbf{A}^T \quad (459)$$

Symmetric matrices have many important properties, e.g. that their eigenvalues are real and eigenvectors orthogonal.

9.8.2 Skew-symmetric/Antisymmetric

The *antisymmetric* matrix is also known as the *skew* symmetric matrix. It has the following property from which it is defined

$$\mathbf{A} = -\mathbf{A}^T \quad (460)$$

Hereby, it can be seen that the antisymmetric matrices always have a zero diagonal. The $n \times n$ antisymmetric matrices also have the following properties.

$$\det(\mathbf{A}^T) = \det(-\mathbf{A}) = (-1)^n \det(\mathbf{A}) \quad (461)$$

$$-\det(\mathbf{A}) = \det(-\mathbf{A}) = 0, \quad \text{if } n \text{ is odd} \quad (462)$$

The eigenvalues of an antisymmetric matrix are placed on the imaginary axis and the eigenvectors are unitary.

9.8.3 Decomposition

A square matrix \mathbf{A} can always be written as a sum of a symmetric \mathbf{A}_+ and an antisymmetric matrix \mathbf{A}_-

$$\mathbf{A} = \mathbf{A}_+ + \mathbf{A}_- \quad (463)$$

Such a decomposition could e.g. be

$$\mathbf{A} = \frac{\mathbf{A} + \mathbf{A}^T}{2} + \frac{\mathbf{A} - \mathbf{A}^T}{2} = \mathbf{A}_+ + \mathbf{A}_- \quad (464)$$

9.9 Toeplitz Matrices

A Toeplitz matrix \mathbf{T} is a matrix where the elements of each diagonal is the same. In the $n \times n$ square case, it has the following structure:

$$\mathbf{T} = \begin{bmatrix} t_{11} & t_{12} & \cdots & t_{1n} \\ t_{21} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & t_{12} \\ t_{n1} & \cdots & t_{21} & t_{11} \end{bmatrix} = \begin{bmatrix} t_0 & t_1 & \cdots & t_{n-1} \\ t_{-1} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & t_1 \\ t_{-(n-1)} & \cdots & t_{-1} & t_0 \end{bmatrix} \quad (465)$$

A Toeplitz matrix is *persymmetric*. If a matrix is persymmetric (or orthosymmetric), it means that the matrix is symmetric about its northeast-southwest diagonal (anti-diagonal) [12]. Persymmetric matrices is a larger class of matrices, since a persymmetric matrix not necessarily has a Toeplitz structure. There

are some special cases of Toeplitz matrices. The *symmetric* Toeplitz matrix is given by:

$$\mathbf{T} = \begin{bmatrix} t_0 & t_1 & \cdots & t_{n-1} \\ t_1 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & t_1 \\ t_{n-1} & \cdots & t_1 & t_0 \end{bmatrix} \quad (466)$$

The circular Toeplitz matrix:

$$\mathbf{T}_C = \begin{bmatrix} t_0 & t_1 & \cdots & t_{n-1} \\ t_{n-1} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & t_1 \\ t_1 & \cdots & t_{n-1} & t_0 \end{bmatrix} \quad (467)$$

The upper triangular Toeplitz matrix:

$$\mathbf{T}_U = \begin{bmatrix} t_0 & t_1 & \cdots & t_{n-1} \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & t_1 \\ 0 & \cdots & 0 & t_0 \end{bmatrix}, \quad (468)$$

and the lower triangular Toeplitz matrix:

$$\mathbf{T}_L = \begin{bmatrix} t_0 & 0 & \cdots & 0 \\ t_{-1} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ t_{-(n-1)} & \cdots & t_{-1} & t_0 \end{bmatrix} \quad (469)$$

9.9.1 Properties of Toeplitz Matrices

The Toeplitz matrix has some computational advantages. The addition of two Toeplitz matrices can be done with $\mathcal{O}(n)$ flops, multiplication of two Toeplitz matrices can be done in $\mathcal{O}(n \ln n)$ flops. Toeplitz equation systems can be solved in $\mathcal{O}(n^2)$ flops. The inverse of a positive definite Toeplitz matrix can be found in $\mathcal{O}(n^2)$ flops too. The inverse of a Toeplitz matrix is persymmetric. The product of two lower triangular Toeplitz matrices is a Toeplitz matrix. More information on Toeplitz matrices and circulant matrices can be found in [13, 7].

9.10 Transition matrices

A square matrix \mathbf{P} is a transition matrix, also known as stochastic matrix or probability matrix, if

$$0 \leq (\mathbf{P})_{ij} \leq 1, \quad \sum_j (\mathbf{P})_{ij} = 1$$

The transition matrix usually describes the probability of moving from state i to j in one step and is closely related to markov processes. Transition matrices

have the following properties

$$\text{Prob}[i \rightarrow j \text{ in 1 step}] = (\mathbf{P})_{ij} \quad (470)$$

$$\text{Prob}[i \rightarrow j \text{ in 2 steps}] = (\mathbf{P}^2)_{ij} \quad (471)$$

$$\text{Prob}[i \rightarrow j \text{ in } k \text{ steps}] = (\mathbf{P}^k)_{ij} \quad (472)$$

$$\text{If all rows are identical} \Rightarrow \mathbf{P}^n = \mathbf{P} \quad (473)$$

$$\boldsymbol{\alpha} \mathbf{P} = \boldsymbol{\alpha}, \quad \boldsymbol{\alpha} \text{ is called invariant} \quad (474)$$

where $\boldsymbol{\alpha}$ is a so-called stationary probability vector, i.e., $0 \leq \alpha_i \leq 1$ and $\sum_i \alpha_i = 1$.

9.11 Units, Permutation and Shift

9.11.1 Unit vector

Let $\mathbf{e}_i \in \mathbb{R}^{n \times 1}$ be the i th unit vector, i.e. the vector which is zero in all entries except the i th at which it is 1.

9.11.2 Rows and Columns

$$i.\text{th row of } \mathbf{A} = \mathbf{e}_i^T \mathbf{A} \quad (475)$$

$$j.\text{th column of } \mathbf{A} = \mathbf{A} \mathbf{e}_j \quad (476)$$

9.11.3 Permutations

Let \mathbf{P} be some permutation matrix, e.g.

$$\mathbf{P} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = [\mathbf{e}_2 \quad \mathbf{e}_1 \quad \mathbf{e}_3] = \begin{bmatrix} \mathbf{e}_2^T \\ \mathbf{e}_1^T \\ \mathbf{e}_3^T \end{bmatrix} \mathbf{A} \quad (477)$$

For permutation matrices it holds that

$$\mathbf{P} \mathbf{P}^T = \mathbf{I} \quad (478)$$

and that

$$\mathbf{A} \mathbf{P} = [\mathbf{A} \mathbf{e}_2 \quad \mathbf{A} \mathbf{e}_1 \quad \mathbf{A} \mathbf{e}_3] \quad \mathbf{P} \mathbf{A} = \begin{bmatrix} \mathbf{e}_2^T \mathbf{A} \\ \mathbf{e}_1^T \mathbf{A} \\ \mathbf{e}_3^T \mathbf{A} \end{bmatrix} \quad (479)$$

That is, the first is a matrix which has columns of \mathbf{A} but in permuted sequence and the second is a matrix which has the rows of \mathbf{A} but in the permuted sequence.

9.11.4 Translation, Shift or Lag Operators

Let \mathbf{L} denote the lag (or 'translation' or 'shift') operator defined on a 4×4 example by

$$\mathbf{L} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad (480)$$

i.e. a matrix of zeros with one on the sub-diagonal, $(\mathbf{L})_{ij} = \delta_{i,j+1}$. With some signal x_t for $t = 1, \dots, N$, the n .th power of the lag operator shifts the indices, i.e.

$$(\mathbf{L}^n \mathbf{x})_t = \begin{cases} 0 & \text{for } t = 1, \dots, n \\ x_{t-n} & \text{for } t = n+1, \dots, N \end{cases} \quad (481)$$

A related but slightly different matrix is the 'recurrent shifted' operator defined on a 4x4 example by

$$\hat{\mathbf{L}} = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad (482)$$

i.e. a matrix defined by $(\hat{\mathbf{L}})_{ij} = \delta_{i,j+1} + \delta_{i,1}\delta_{j,dim(\mathbf{L})}$. On a signal \mathbf{x} it has the effect

$$(\hat{\mathbf{L}}^n \mathbf{x})_t = x_{t'}, \quad t' = [(t-n) \bmod N] + 1 \quad (483)$$

That is, $\hat{\mathbf{L}}$ is like the shift operator \mathbf{L} except that it 'wraps' the signal as if it was periodic and shifted (substituting the zeros with the rear end of the signal). Note that $\hat{\mathbf{L}}$ is invertible and orthogonal, i.e.

$$\hat{\mathbf{L}}^{-1} = \hat{\mathbf{L}}^T \quad (484)$$

9.12 Vandermonde Matrices

A Vandermonde matrix has the form [15]

$$\mathbf{V} = \begin{bmatrix} 1 & v_1 & v_1^2 & \dots & v_1^{n-1} \\ 1 & v_2 & v_2^2 & \dots & v_2^{n-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & v_n & v_n^2 & \dots & v_n^{n-1} \end{bmatrix}. \quad (485)$$

The transpose of \mathbf{V} is also said to a Vandermonde matrix. The determinant is given by [29]

$$\det \mathbf{V} = \prod_{i>j} (v_i - v_j) \quad (486)$$

10 Functions and Operators

10.1 Functions and Series

10.1.1 Finite Series

$$(\mathbf{X}^n - \mathbf{I})(\mathbf{X} - \mathbf{I})^{-1} = \mathbf{I} + \mathbf{X} + \mathbf{X}^2 + \dots + \mathbf{X}^{n-1} \quad (487)$$

10.1.2 Taylor Expansion of Scalar Function

Consider some scalar function $f(\mathbf{x})$ which takes the vector \mathbf{x} as an argument. This we can Taylor expand around \mathbf{x}_0

$$f(\mathbf{x}) \cong f(\mathbf{x}_0) + \mathbf{g}(\mathbf{x}_0)^T(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{H}(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \quad (488)$$

where

$$\mathbf{g}(\mathbf{x}_0) = \left. \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_0} \quad \mathbf{H}(\mathbf{x}_0) = \left. \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} \right|_{\mathbf{x}_0}$$

10.1.3 Matrix Functions by Infinite Series

As for analytical functions in one dimension, one can define a matrix function for square matrices \mathbf{X} by an infinite series

$$\mathbf{f}(\mathbf{X}) = \sum_{n=0}^{\infty} c_n \mathbf{X}^n \quad (489)$$

assuming the limit exists and is finite. If the coefficients c_n fulfils $\sum_n c_n x^n < \infty$, then one can prove that the above series exists and is finite, see [1]. Thus for any analytical function $f(x)$ there exists a corresponding matrix function $\mathbf{f}(\mathbf{x})$ constructed by the Taylor expansion. Using this one can prove the following results:

1) A matrix \mathbf{A} is a zero of its own characteristic polynomial [1]:

$$p(\lambda) = \det(\mathbf{I}\lambda - \mathbf{A}) = \sum_n c_n \lambda^n \quad \Rightarrow \quad p(\mathbf{A}) = \mathbf{0} \quad (490)$$

2) If \mathbf{A} is square it holds that [1]

$$\mathbf{A} = \mathbf{U}\mathbf{B}\mathbf{U}^{-1} \quad \Rightarrow \quad \mathbf{f}(\mathbf{A}) = \mathbf{U}\mathbf{f}(\mathbf{B})\mathbf{U}^{-1} \quad (491)$$

3) A useful fact when using power series is that

$$\mathbf{A}^n \rightarrow \mathbf{0} \text{ for } n \rightarrow \infty \quad \text{if} \quad |\mathbf{A}| < 1 \quad (492)$$

10.1.4 Identity and commutations

It holds for an analytical matrix function $\mathbf{f}(\mathbf{X})$ that

$$\mathbf{f}(\mathbf{A}\mathbf{B})\mathbf{A} = \mathbf{A}\mathbf{f}(\mathbf{B}\mathbf{A}) \quad (493)$$

see B.1.2 for a proof.

10.1.5 Exponential Matrix Function

In analogy to the ordinary scalar exponential function, one can define exponential and logarithmic matrix functions:

$$e^{\mathbf{A}} \equiv \sum_{n=0}^{\infty} \frac{1}{n!} \mathbf{A}^n = \mathbf{I} + \mathbf{A} + \frac{1}{2} \mathbf{A}^2 + \dots \quad (494)$$

$$e^{-\mathbf{A}} \equiv \sum_{n=0}^{\infty} \frac{1}{n!} (-1)^n \mathbf{A}^n = \mathbf{I} - \mathbf{A} + \frac{1}{2} \mathbf{A}^2 - \dots \quad (495)$$

$$e^{t\mathbf{A}} \equiv \sum_{n=0}^{\infty} \frac{1}{n!} (t\mathbf{A})^n = \mathbf{I} + t\mathbf{A} + \frac{1}{2} t^2 \mathbf{A}^2 + \dots \quad (496)$$

$$\ln(\mathbf{I} + \mathbf{A}) \equiv \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} \mathbf{A}^n = \mathbf{A} - \frac{1}{2} \mathbf{A}^2 + \frac{1}{3} \mathbf{A}^3 - \dots \quad (497)$$

Some of the properties of the exponential function are [1]

$$e^{\mathbf{A}} e^{\mathbf{B}} = e^{\mathbf{A}+\mathbf{B}} \quad \text{if} \quad \mathbf{AB} = \mathbf{BA} \quad (498)$$

$$(e^{\mathbf{A}})^{-1} = e^{-\mathbf{A}} \quad (499)$$

$$\frac{d}{dt} e^{t\mathbf{A}} = \mathbf{A} e^{t\mathbf{A}} = e^{t\mathbf{A}} \mathbf{A}, \quad t \in \mathbb{R} \quad (500)$$

$$\frac{d}{dt} \text{Tr}(e^{t\mathbf{A}}) = \text{Tr}(\mathbf{A} e^{t\mathbf{A}}) \quad (501)$$

$$\det(e^{\mathbf{A}}) = e^{\text{Tr}(\mathbf{A})} \quad (502)$$

10.1.6 Trigonometric Functions

$$\sin(\mathbf{A}) \equiv \sum_{n=0}^{\infty} \frac{(-1)^n \mathbf{A}^{2n+1}}{(2n+1)!} = \mathbf{A} - \frac{1}{3!} \mathbf{A}^3 + \frac{1}{5!} \mathbf{A}^5 - \dots \quad (503)$$

$$\cos(\mathbf{A}) \equiv \sum_{n=0}^{\infty} \frac{(-1)^n \mathbf{A}^{2n}}{(2n)!} = \mathbf{I} - \frac{1}{2!} \mathbf{A}^2 + \frac{1}{4!} \mathbf{A}^4 - \dots \quad (504)$$

10.2 Kronecker and Vec Operator

10.2.1 The Kronecker Product

The Kronecker product of an $m \times n$ matrix \mathbf{A} and an $r \times q$ matrix \mathbf{B} , is an $mr \times nq$ matrix, $\mathbf{A} \otimes \mathbf{B}$ defined as

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} A_{11}\mathbf{B} & A_{12}\mathbf{B} & \dots & A_{1n}\mathbf{B} \\ A_{21}\mathbf{B} & A_{22}\mathbf{B} & \dots & A_{2n}\mathbf{B} \\ \vdots & & & \vdots \\ A_{m1}\mathbf{B} & A_{m2}\mathbf{B} & \dots & A_{mn}\mathbf{B} \end{bmatrix} \quad (505)$$

The Kronecker product has the following properties (see [19])

$$\mathbf{A} \otimes (\mathbf{B} + \mathbf{C}) = \mathbf{A} \otimes \mathbf{B} + \mathbf{A} \otimes \mathbf{C} \quad (506)$$

$$\mathbf{A} \otimes \mathbf{B} \neq \mathbf{B} \otimes \mathbf{A} \quad \text{in general} \quad (507)$$

$$\mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}) = (\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} \quad (508)$$

$$(\alpha_A \mathbf{A} \otimes \alpha_B \mathbf{B}) = \alpha_A \alpha_B (\mathbf{A} \otimes \mathbf{B}) \quad (509)$$

$$(\mathbf{A} \otimes \mathbf{B})^T = \mathbf{A}^T \otimes \mathbf{B}^T \quad (510)$$

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD} \quad (511)$$

$$(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1} \quad (512)$$

$$(\mathbf{A} \otimes \mathbf{B})^+ = \mathbf{A}^+ \otimes \mathbf{B}^+ \quad (513)$$

$$\text{rank}(\mathbf{A} \otimes \mathbf{B}) = \text{rank}(\mathbf{A})\text{rank}(\mathbf{B}) \quad (514)$$

$$\text{Tr}(\mathbf{A} \otimes \mathbf{B}) = \text{Tr}(\mathbf{A})\text{Tr}(\mathbf{B}) = \text{Tr}(\mathbf{\Lambda}_A \otimes \mathbf{\Lambda}_B) \quad (515)$$

$$\det(\mathbf{A} \otimes \mathbf{B}) = \det(\mathbf{A})^{\text{rank}(\mathbf{B})} \det(\mathbf{B})^{\text{rank}(\mathbf{A})} \quad (516)$$

$$\{\text{eig}(\mathbf{A} \otimes \mathbf{B})\} = \{\text{eig}(\mathbf{B} \otimes \mathbf{A})\} \quad \text{if } \mathbf{A}, \mathbf{B} \text{ are square} \quad (517)$$

$$\{\text{eig}(\mathbf{A} \otimes \mathbf{B})\} = \{\text{eig}(\mathbf{A})\text{eig}(\mathbf{B})^T\} \quad (518)$$

if \mathbf{A}, \mathbf{B} are symmetric and square

$$\text{eig}(\mathbf{A} \otimes \mathbf{B}) = \text{eig}(\mathbf{A}) \otimes \text{eig}(\mathbf{B}) \quad (519)$$

Where $\{\lambda_i\}$ denotes the set of values λ_i , that is, the values in no particular order or structure, and $\mathbf{\Lambda}_A$ denotes the diagonal matrix with the eigenvalues of \mathbf{A} .

10.2.2 The Vec Operator

The vec-operator applied on a matrix \mathbf{A} stacks the columns into a vector, i.e. for a 2×2 matrix

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \text{vec}(\mathbf{A}) = \begin{bmatrix} A_{11} \\ A_{21} \\ A_{12} \\ A_{22} \end{bmatrix}$$

Properties of the vec-operator include (see [19])

$$\text{vec}(\mathbf{AXB}) = (\mathbf{B}^T \otimes \mathbf{A})\text{vec}(\mathbf{X}) \quad (520)$$

$$\text{Tr}(\mathbf{A}^T \mathbf{B}) = \text{vec}(\mathbf{A})^T \text{vec}(\mathbf{B}) \quad (521)$$

$$\text{vec}(\mathbf{A} + \mathbf{B}) = \text{vec}(\mathbf{A}) + \text{vec}(\mathbf{B}) \quad (522)$$

$$\text{vec}(\alpha \mathbf{A}) = \alpha \cdot \text{vec}(\mathbf{A}) \quad (523)$$

$$\mathbf{a}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{c} = \text{vec}(\mathbf{X})^T (\mathbf{B} \otimes \mathbf{ca}^T) \text{vec}(\mathbf{X}) \quad (524)$$

See B.1.1 for a proof for Eq. 524.

10.3 Vector Norms

10.3.1 Examples

$$\|\mathbf{x}\|_1 = \sum_i |x_i| \quad (525)$$

$$\|\mathbf{x}\|_2^2 = \mathbf{x}^H \mathbf{x} \quad (526)$$

$$\|\mathbf{x}\|_p = \left[\sum_i |x_i|^p \right]^{1/p} \quad (527)$$

$$\|\mathbf{x}\|_\infty = \max_i |x_i| \quad (528)$$

Further reading in e.g. [12, p. 52]

10.4 Matrix Norms

10.4.1 Definitions

A matrix norm is a mapping which fulfils

$$\|\mathbf{A}\| \geq 0 \quad (529)$$

$$\|\mathbf{A}\| = 0 \Leftrightarrow \mathbf{A} = \mathbf{0} \quad (530)$$

$$\|c\mathbf{A}\| = |c| \|\mathbf{A}\|, \quad c \in \mathbb{R} \quad (531)$$

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\| \quad (532)$$

10.4.2 Induced Norm or Operator Norm

An induced norm is a matrix norm induced by a vector norm by the following

$$\|\mathbf{A}\| = \sup\{\|\mathbf{Ax}\| \mid \|\mathbf{x}\| = 1\} \quad (533)$$

where $\|\cdot\|$ on the left side is the induced matrix norm, while $\|\cdot\|$ on the right side denotes the vector norm. For induced norms it holds that

$$\|\mathbf{I}\| = 1 \quad (534)$$

$$\|\mathbf{Ax}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|, \quad \text{for all } \mathbf{A}, \mathbf{x} \quad (535)$$

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|, \quad \text{for all } \mathbf{A}, \mathbf{B} \quad (536)$$

10.4.3 Examples

$$\|\mathbf{A}\|_1 = \max_j \sum_i |A_{ij}| \quad (537)$$

$$\|\mathbf{A}\|_2 = \sqrt{\max \text{eig}(\mathbf{A}^H \mathbf{A})} \quad (538)$$

$$\|\mathbf{A}\|_p = \left(\max_{\|\mathbf{x}\|_p=1} \|\mathbf{Ax}\|_p \right)^{1/p} \quad (539)$$

$$\|\mathbf{A}\|_\infty = \max_i \sum_j |A_{ij}| \quad (540)$$

$$\|\mathbf{A}\|_F = \sqrt{\sum_{ij} |A_{ij}|^2} = \sqrt{\text{Tr}(\mathbf{AA}^H)} \quad (\text{Frobenius}) \quad (541)$$

$$\|\mathbf{A}\|_{\max} = \max_{ij} |A_{ij}| \quad (542)$$

$$\|\mathbf{A}\|_{\text{KF}} = \|\text{sing}(\mathbf{A})\|_1 \quad (\text{Ky Fan}) \quad (543)$$

where $\text{sing}(\mathbf{A})$ is the vector of singular values of the matrix \mathbf{A} .

10.4.4 Inequalities

E. H. Rasmussen has in yet unpublished material derived and collected the following inequalities. They are collected in a table as below, assuming \mathbf{A} is an $m \times n$, and $d = \text{rank}(\mathbf{A})$

	$\ \mathbf{A}\ _{\max}$	$\ \mathbf{A}\ _1$	$\ \mathbf{A}\ _{\infty}$	$\ \mathbf{A}\ _2$	$\ \mathbf{A}\ _{\text{F}}$	$\ \mathbf{A}\ _{\text{KF}}$
$\ \mathbf{A}\ _{\max}$		1	1	1	1	1
$\ \mathbf{A}\ _1$	m		m	\sqrt{m}	\sqrt{m}	\sqrt{m}
$\ \mathbf{A}\ _{\infty}$	n	n		\sqrt{n}	\sqrt{n}	\sqrt{n}
$\ \mathbf{A}\ _2$	\sqrt{mn}	\sqrt{n}	\sqrt{m}		1	1
$\ \mathbf{A}\ _{\text{F}}$	\sqrt{mn}	\sqrt{n}	\sqrt{m}	\sqrt{d}		1
$\ \mathbf{A}\ _{\text{KF}}$	\sqrt{mnd}	\sqrt{nd}	\sqrt{md}	d	\sqrt{d}	

which are to be read as, e.g.

$$\|\mathbf{A}\|_2 \leq \sqrt{m} \cdot \|\mathbf{A}\|_{\infty} \quad (544)$$

10.4.5 Condition Number

The 2-norm of \mathbf{A} equals $\sqrt{(\max(\text{eig}(\mathbf{A}^T \mathbf{A}))})}$ [12, p.57]. For a symmetric, positive definite matrix, this reduces to $\max(\text{eig}(\mathbf{A}))$. The condition number based on the 2-norm thus reduces to

$$\|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \max(\text{eig}(\mathbf{A})) \max(\text{eig}(\mathbf{A}^{-1})) = \frac{\max(\text{eig}(\mathbf{A}))}{\min(\text{eig}(\mathbf{A}))}. \quad (545)$$

10.5 Rank

10.5.1 Sylvester's Inequality

If \mathbf{A} is $m \times n$ and \mathbf{B} is $n \times r$, then

$$\text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B}) - n \leq \text{rank}(\mathbf{AB}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\} \quad (546)$$

10.6 Integral Involving Dirac Delta Functions

Assuming \mathbf{A} to be square, then

$$\int p(\mathbf{s}) \delta(\mathbf{x} - \mathbf{As}) d\mathbf{s} = \frac{1}{\det(\mathbf{A})} p(\mathbf{A}^{-1} \mathbf{x}) \quad (547)$$

Assuming \mathbf{A} to be "underdetermined", i.e. "tall", then

$$\int p(\mathbf{s}) \delta(\mathbf{x} - \mathbf{As}) d\mathbf{s} = \begin{cases} \frac{1}{\sqrt{\det(\mathbf{A}^T \mathbf{A})}} p(\mathbf{A}^+ \mathbf{x}) & \text{if } \mathbf{x} = \mathbf{AA}^+ \mathbf{x} \\ 0 & \text{elsewhere} \end{cases} \quad (548)$$

See [9].

10.7 Miscellaneous

For any \mathbf{A} it holds that

$$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^T) = \text{rank}(\mathbf{A}\mathbf{A}^T) = \text{rank}(\mathbf{A}^T\mathbf{A}) \quad (549)$$

It holds that

$$\mathbf{A} \text{ is positive definite} \Leftrightarrow \exists \mathbf{B} \text{ invertible, such that } \mathbf{A} = \mathbf{B}\mathbf{B}^T \quad (550)$$

A One-dimensional Results

A.1 Gaussian

A.1.1 Density

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (551)$$

A.1.2 Normalization

$$\int e^{-\frac{(s-\mu)^2}{2\sigma^2}} ds = \sqrt{2\pi\sigma^2} \quad (552)$$

$$\int e^{-(ax^2+bx+c)} dx = \sqrt{\frac{\pi}{a}} \exp\left[\frac{b^2-4ac}{4a}\right] \quad (553)$$

$$\int e^{c_2x^2+c_1x+c_0} dx = \sqrt{\frac{\pi}{-c_2}} \exp\left[\frac{c_1^2-4c_2c_0}{-4c_2}\right] \quad (554)$$

A.1.3 Derivatives

$$\frac{\partial p(x)}{\partial \mu} = p(x) \frac{(x-\mu)}{\sigma^2} \quad (555)$$

$$\frac{\partial \ln p(x)}{\partial \mu} = \frac{(x-\mu)}{\sigma^2} \quad (556)$$

$$\frac{\partial p(x)}{\partial \sigma} = p(x) \frac{1}{\sigma} \left[\frac{(x-\mu)^2}{\sigma^2} - 1 \right] \quad (557)$$

$$\frac{\partial \ln p(x)}{\partial \sigma} = \frac{1}{\sigma} \left[\frac{(x-\mu)^2}{\sigma^2} - 1 \right] \quad (558)$$

A.1.4 Completing the Squares

$$c_2x^2 + c_1x + c_0 = -a(x-b)^2 + w$$

$$-a = c_2 \quad b = \frac{1}{2} \frac{c_1}{c_2} \quad w = \frac{1}{4} \frac{c_1^2}{c_2} + c_0$$

or

$$c_2x^2 + c_1x + c_0 = -\frac{1}{2\sigma^2}(x-\mu)^2 + d$$

$$\mu = \frac{-c_1}{2c_2} \quad \sigma^2 = \frac{-1}{2c_2} \quad d = c_0 - \frac{c_1^2}{4c_2}$$

A.1.5 Moments

If the density is expressed by

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(s-\mu)^2}{2\sigma^2}\right] \quad \text{or} \quad p(x) = C \exp(c_2x^2 + c_1x) \quad (559)$$

then the first few basic moments are

$$\begin{aligned}
 \langle x \rangle &= \mu &= \frac{-c_1}{2c_2} \\
 \langle x^2 \rangle &= \sigma^2 + \mu^2 &= \frac{-1}{2c_2} + \left(\frac{-c_1}{2c_2} \right)^2 \\
 \langle x^3 \rangle &= 3\sigma^2\mu + \mu^3 &= \frac{c_1}{(2c_2)^2} \left[3 - \frac{c_1^2}{2c_2} \right] \\
 \langle x^4 \rangle &= \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4 &= \left(\frac{c_1}{2c_2} \right)^4 + 6 \left(\frac{c_1}{2c_2} \right)^2 \left(\frac{-1}{2c_2} \right) + 3 \left(\frac{1}{2c_2} \right)^2
 \end{aligned}$$

and the central moments are

$$\begin{aligned}
 \langle (x - \mu) \rangle &= 0 &= 0 \\
 \langle (x - \mu)^2 \rangle &= \sigma^2 &= \left[\frac{-1}{2c_2} \right] \\
 \langle (x - \mu)^3 \rangle &= 0 &= 0 \\
 \langle (x - \mu)^4 \rangle &= 3\sigma^4 &= 3 \left[\frac{1}{2c_2} \right]^2
 \end{aligned}$$

A kind of pseudo-moments (un-normalized integrals) can easily be derived as

$$\int \exp(c_2 x^2 + c_1 x) x^n dx = Z \langle x^n \rangle = \sqrt{\frac{\pi}{-c_2}} \exp \left[\frac{c_1^2}{-4c_2} \right] \langle x^n \rangle \quad (560)$$

From the un-centralized moments one can derive other entities like

$$\begin{aligned}
 \langle x^2 \rangle - \langle x \rangle^2 &= \sigma^2 &= \frac{-1}{2c_2} \\
 \langle x^3 \rangle - \langle x^2 \rangle \langle x \rangle &= 2\sigma^2\mu &= \frac{2c_1}{(2c_2)^2} \\
 \langle x^4 \rangle - \langle x^2 \rangle^2 &= 2\sigma^4 + 4\mu^2\sigma^2 &= \frac{2}{(2c_2)^2} \left[1 - 4 \frac{c_1^2}{2c_2} \right]
 \end{aligned}$$

A.2 One Dimensional Mixture of Gaussians

A.2.1 Density and Normalization

$$p(s) = \sum_k^K \frac{\rho_k}{\sqrt{2\pi\sigma_k^2}} \exp \left[-\frac{1}{2} \frac{(s - \mu_k)^2}{\sigma_k^2} \right] \quad (561)$$

A.2.2 Moments

A useful fact of MoG, is that

$$\langle x^n \rangle = \sum_k \rho_k \langle x^n \rangle_k \quad (562)$$

where $\langle \cdot \rangle_k$ denotes average with respect to the k .th component. We can calculate the first four moments from the densities

$$p(x) = \sum_k \rho_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp \left[-\frac{1}{2} \frac{(x - \mu_k)^2}{\sigma_k^2} \right] \quad (563)$$

$$p(x) = \sum_k \rho_k C_k \exp [c_{k2}x^2 + c_{k1}x] \quad (564)$$

as

$$\begin{aligned}
 \langle x \rangle &= \sum_k \rho_k \mu_k &= \sum_k \rho_k \left[\frac{-c_{k1}}{2c_{k2}} \right] \\
 \langle x^2 \rangle &= \sum_k \rho_k (\sigma_k^2 + \mu_k^2) &= \sum_k \rho_k \left[\frac{-1}{2c_{k2}} + \left(\frac{-c_{k1}}{2c_{k2}} \right)^2 \right] \\
 \langle x^3 \rangle &= \sum_k \rho_k (3\sigma_k^2 \mu_k + \mu_k^3) &= \sum_k \rho_k \left[\frac{c_{k1}}{(2c_{k2})^2} \left[3 - \frac{c_{k1}^2}{2c_{k2}} \right] \right] \\
 \langle x^4 \rangle &= \sum_k \rho_k (\mu_k^4 + 6\mu_k^2 \sigma_k^2 + 3\sigma_k^4) &= \sum_k \rho_k \left[\left(\frac{1}{2c_{k2}} \right)^2 \left[\left(\frac{c_{k1}}{2c_{k2}} \right)^2 - 6 \frac{c_{k1}^2}{2c_{k2}} + 3 \right] \right]
 \end{aligned}$$

If all the gaussians are centered, i.e. $\mu_k = 0$ for all k , then

$$\begin{aligned}
 \langle x \rangle &= 0 &= 0 \\
 \langle x^2 \rangle &= \sum_k \rho_k \sigma_k^2 &= \sum_k \rho_k \left[\frac{-1}{2c_{k2}} \right] \\
 \langle x^3 \rangle &= 0 &= 0 \\
 \langle x^4 \rangle &= \sum_k \rho_k 3\sigma_k^4 &= \sum_k \rho_k 3 \left[\frac{-1}{2c_{k2}} \right]^2
 \end{aligned}$$

From the un-centralized moments one can derive other entities like

$$\begin{aligned}
 \langle x^2 \rangle - \langle x \rangle^2 &= \sum_{k,k'} \rho_k \rho_{k'} [\mu_k^2 + \sigma_k^2 - \mu_k \mu_{k'}] \\
 \langle x^3 \rangle - \langle x^2 \rangle \langle x \rangle &= \sum_{k,k'} \rho_k \rho_{k'} [3\sigma_k^2 \mu_k + \mu_k^3 - (\sigma_k^2 + \mu_k^2) \mu_{k'}] \\
 \langle x^4 \rangle - \langle x^2 \rangle^2 &= \sum_{k,k'} \rho_k \rho_{k'} [\mu_k^4 + 6\mu_k^2 \sigma_k^2 + 3\sigma_k^4 - (\sigma_k^2 + \mu_k^2)(\sigma_{k'}^2 + \mu_{k'}^2)]
 \end{aligned}$$

A.2.3 Derivatives

Defining $p(s) = \sum_k \rho_k \mathcal{N}_s(\mu_k, \sigma_k^2)$ we get for a parameter θ_j of the j .th component

$$\frac{\partial \ln p(s)}{\partial \theta_j} = \frac{\rho_j \mathcal{N}_s(\mu_j, \sigma_j^2)}{\sum_k \rho_k \mathcal{N}_s(\mu_k, \sigma_k^2)} \frac{\partial \ln(\rho_j \mathcal{N}_s(\mu_j, \sigma_j^2))}{\partial \theta_j} \quad (565)$$

that is,

$$\frac{\partial \ln p(s)}{\partial \rho_j} = \frac{\rho_j \mathcal{N}_s(\mu_j, \sigma_j^2)}{\sum_k \rho_k \mathcal{N}_s(\mu_k, \sigma_k^2)} \frac{1}{\rho_j} \quad (566)$$

$$\frac{\partial \ln p(s)}{\partial \mu_j} = \frac{\rho_j \mathcal{N}_s(\mu_j, \sigma_j^2)}{\sum_k \rho_k \mathcal{N}_s(\mu_k, \sigma_k^2)} \frac{(s - \mu_j)}{\sigma_j^2} \quad (567)$$

$$\frac{\partial \ln p(s)}{\partial \sigma_j} = \frac{\rho_j \mathcal{N}_s(\mu_j, \sigma_j^2)}{\sum_k \rho_k \mathcal{N}_s(\mu_k, \sigma_k^2)} \frac{1}{\sigma_j} \left[\frac{(s - \mu_j)^2}{\sigma_j^2} - 1 \right] \quad (568)$$

Note that ρ_k must be constrained to be proper ratios. Defining the ratios by $\rho_j = e^{r_j} / \sum_k e^{r_k}$, we obtain

$$\frac{\partial \ln p(s)}{\partial r_j} = \sum_l \frac{\partial \ln p(s)}{\partial \rho_l} \frac{\partial \rho_l}{\partial r_j} \quad \text{where} \quad \frac{\partial \rho_l}{\partial r_j} = \rho_l (\delta_{lj} - \rho_j) \quad (569)$$

B Proofs and Details

B.1 Misc Proofs

B.1.1 Proof of Equation 524

The following proof is work of Florian Roemer. Note the the vectors and matrices below can be complex and the notation \mathbf{X}^H is used for transpose and conjugated, while \mathbf{X}^T is *only* transpose of the complex matrix.

Define the row vector $\mathbf{y} = \mathbf{a}^H \mathbf{X} \mathbf{B}$ and the column vector $\mathbf{z} = \mathbf{X}^H \mathbf{c}$. Then

$$\mathbf{a}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{c} = \mathbf{y} \mathbf{z} = \mathbf{z}^T \mathbf{y}^T$$

Note that \mathbf{y} can be rewritten as $\text{vec}(\mathbf{y})^T$ which is the same as

$$\text{vec}(\text{conj}(\mathbf{y}))^H = \text{vec}(\mathbf{a}^T \text{conj}(\mathbf{X}) \text{conj}(\mathbf{B}))^H$$

where "conj" means complex conjugated. Applying the vec rule for linear forms Eq 520, we get

$$\mathbf{y} = (\mathbf{B}^H \otimes \mathbf{a}^T \text{vec}(\text{conj}(\mathbf{X})))^H = \text{vec}(\mathbf{X})^T (\mathbf{B} \otimes \text{conj}(\mathbf{a}))$$

where we have also used the rule for transpose of Kronecker products. For \mathbf{y}^T this yields $(\mathbf{B}^T \otimes \mathbf{a}^H) \text{vec}(\mathbf{X})$. Similarly we can rewrite \mathbf{z} which is the same as $\text{vec}(\mathbf{z}^T) = \text{vec}(\mathbf{c}^T \text{conj}(\mathbf{X}))$. Applying again Eq 520, we get

$$\mathbf{z} = (\mathbf{I} \otimes \mathbf{c}^T) \text{vec}(\text{conj}(\mathbf{X}))$$

where \mathbf{I} is the identity matrix. For \mathbf{z}^T we obtain $\text{vec}(\mathbf{X})(\mathbf{I} \otimes \mathbf{c})$. Finally, the original expression is $\mathbf{z}^T \mathbf{y}^T$ which now takes the form

$$\text{vec}(\mathbf{X})^H (\mathbf{I} \otimes \mathbf{c}) (\mathbf{B}^T \otimes \mathbf{a}^H) \text{vec}(\mathbf{X})$$

the final step is to apply the rule for products of Kronecker products and by that combine the Kronecker products. This gives

$$\text{vec}(\mathbf{X})^H (\mathbf{B}^T \otimes \mathbf{c} \mathbf{a}^H) \text{vec}(\mathbf{X})$$

which is the desired result.

B.1.2 Proof of Equation 493

For any analytical function $\mathbf{f}(\mathbf{X})$ of a matrix argument \mathbf{X} , it holds that

$$\begin{aligned} \mathbf{f}(\mathbf{A} \mathbf{B}) \mathbf{A} &= \left(\sum_{n=0}^{\infty} c_n (\mathbf{A} \mathbf{B})^n \right) \mathbf{A} \\ &= \sum_{n=0}^{\infty} c_n (\mathbf{A} \mathbf{B})^n \mathbf{A} \\ &= \sum_{n=0}^{\infty} c_n \mathbf{A} (\mathbf{B} \mathbf{A})^n \\ &= \mathbf{A} \sum_{n=0}^{\infty} c_n (\mathbf{B} \mathbf{A})^n \\ &= \mathbf{A} \mathbf{f}(\mathbf{B} \mathbf{A}) \end{aligned}$$

B.1.3 Proof of Equation 91

Essentially we need to calculate

$$\begin{aligned}
\frac{\partial(\mathbf{X}^n)_{kl}}{\partial X_{ij}} &= \frac{\partial}{\partial X_{ij}} \sum_{u_1, \dots, u_{n-1}} X_{k,u_1} X_{u_1,u_2} \dots X_{u_{n-1},l} \\
&= \delta_{k,i} \delta_{u_1,j} X_{u_1,u_2} \dots X_{u_{n-1},l} \\
&\quad + X_{k,u_1} \delta_{u_1,i} \delta_{u_2,j} \dots X_{u_{n-1},l} \\
&\quad \vdots \\
&\quad + X_{k,u_1} X_{u_1,u_2} \dots \delta_{u_{n-1},i} \delta_{l,j} \\
&= \sum_{r=0}^{n-1} (\mathbf{X}^r)_{ki} (\mathbf{X}^{n-1-r})_{jl} \\
&= \sum_{r=0}^{n-1} (\mathbf{X}^r \mathbf{J}^{ij} \mathbf{X}^{n-1-r})_{kl}
\end{aligned}$$

Using the properties of the single entry matrix found in Sec. 9.7.4, the result follows easily.

B.1.4 Details on Eq. 571

$$\begin{aligned}
\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) &= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) \text{Tr}[(\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \partial(\mathbf{X}^H \mathbf{A} \mathbf{X})] \\
&= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) \text{Tr}[(\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} (\partial(\mathbf{X}^H) \mathbf{A} \mathbf{X} + \mathbf{X}^H \partial(\mathbf{A} \mathbf{X}))] \\
&= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) (\text{Tr}[(\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \partial(\mathbf{X}^H) \mathbf{A} \mathbf{X}] \\
&\quad + \text{Tr}[(\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \partial(\mathbf{A} \mathbf{X})]) \\
&= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) (\text{Tr}[\mathbf{A} \mathbf{X} (\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \partial(\mathbf{X}^H)] \\
&\quad + \text{Tr}[(\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \mathbf{A} \partial(\mathbf{X})])
\end{aligned}$$

First, the derivative is found with respect to the real part of \mathbf{X}

$$\begin{aligned}
\frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Re \mathbf{X}} &= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) \left(\frac{\text{Tr}[\mathbf{A} \mathbf{X} (\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \partial(\mathbf{X}^H)]}{\partial \Re \mathbf{X}} \right. \\
&\quad \left. + \frac{\text{Tr}[(\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \mathbf{A} \partial(\mathbf{X})]}{\partial \Re \mathbf{X}} \right) \\
&= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) (\mathbf{A} \mathbf{X} (\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} + ((\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \mathbf{A})^T)
\end{aligned}$$

Through the calculations, (100) and (240) were used. In addition, by use of (241), the derivative is found with respect to the imaginary part of \mathbf{X}

$$\begin{aligned}
i \frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Im \mathbf{X}} &= i \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) \left(\frac{\text{Tr}[\mathbf{A} \mathbf{X} (\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \partial(\mathbf{X}^H)]}{\partial \Im \mathbf{X}} \right. \\
&\quad \left. + \frac{\text{Tr}[(\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \mathbf{A} \partial(\mathbf{X})]}{\partial \Im \mathbf{X}} \right) \\
&= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) (\mathbf{A} \mathbf{X} (\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} - ((\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \mathbf{A})^T)
\end{aligned}$$

Hence, derivative yields

$$\begin{aligned}
\frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \mathbf{X}} &= \frac{1}{2} \left(\frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Re \mathbf{X}} - i \frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Im \mathbf{X}} \right) \\
&= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) ((\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^H \mathbf{A})^T
\end{aligned}$$

and the complex conjugate derivative yields

$$\begin{aligned}\frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \mathbf{X}^*} &= \frac{1}{2} \left(\frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Re \mathbf{X}} + i \frac{\partial \det(\mathbf{X}^H \mathbf{A} \mathbf{X})}{\partial \Im \mathbf{X}} \right) \\ &= \det(\mathbf{X}^H \mathbf{A} \mathbf{X}) \mathbf{A} \mathbf{X} (\mathbf{X}^H \mathbf{A} \mathbf{X})^{-1}\end{aligned}$$

Notice, for real \mathbf{X} , \mathbf{A} , the sum of (249) and (250) is reduced to (54). Similar calculations yield

$$\begin{aligned}\frac{\partial \det(\mathbf{X} \mathbf{A} \mathbf{X}^H)}{\partial \mathbf{X}} &= \frac{1}{2} \left(\frac{\partial \det(\mathbf{X} \mathbf{A} \mathbf{X}^H)}{\partial \Re \mathbf{X}} - i \frac{\partial \det(\mathbf{X} \mathbf{A} \mathbf{X}^H)}{\partial \Im \mathbf{X}} \right) \\ &= \det(\mathbf{X} \mathbf{A} \mathbf{X}^H) (\mathbf{A} \mathbf{X}^H (\mathbf{X} \mathbf{A} \mathbf{X}^H)^{-1})^T\end{aligned}\quad (570)$$

and

$$\begin{aligned}\frac{\partial \det(\mathbf{X} \mathbf{A} \mathbf{X}^H)}{\partial \mathbf{X}^*} &= \frac{1}{2} \left(\frac{\partial \det(\mathbf{X} \mathbf{A} \mathbf{X}^H)}{\partial \Re \mathbf{X}} + i \frac{\partial \det(\mathbf{X} \mathbf{A} \mathbf{X}^H)}{\partial \Im \mathbf{X}} \right) \\ &= \det(\mathbf{X} \mathbf{A} \mathbf{X}^H) (\mathbf{X} \mathbf{A} \mathbf{X}^H)^{-1} \mathbf{X} \mathbf{A}\end{aligned}\quad (571)$$

References

- [1] Karl Gustav Andersson and Lars-Christer Boiers. *Ordinæra differentialekvationer*. Studentlitteratur, 1992.
- [2] Jörn Anemüller, Terrence J. Sejnowski, and Scott Makeig. Complex independent component analysis of frequency-domain electroencephalographic data. *Neural Networks*, 16(9):1311–1323, November 2003.
- [3] S. Barnett. *Matrices. Methods and Applications*. Oxford Applied Mathematics and Computing Science Series. Clarendon Press, 1990.
- [4] Christopher Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- [5] Robert J. Boik. Lecture notes: Statistics 550. Online, April 22 2002. Notes.
- [6] D. H. Brandwood. A complex gradient operator and its application in adaptive array theory. *IEEE Proceedings*, 130(1):11–16, February 1983. PTS. F and H.
- [7] M. Brookes. Matrix Reference Manual, 2004. Website May 20, 2004.
- [8] Contradsen K., *En introduktion til statistik*, IMM lecture notes, 1984.
- [9] Mads Dyrholm. Some matrix results, 2004. Website August 23, 2004.
- [10] Nielsen F. A., *Formula*, Neuro Research Unit and Technical university of Denmark, 2002.
- [11] Gelman A. B., J. S. Carlin, H. S. Stern, D. B. Rubin, *Bayesian Data Analysis*, Chapman and Hall / CRC, 1995.
- [12] Gene H. Golub and Charles F. van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, 3rd edition, 1996.
- [13] Robert M. Gray. Toeplitz and circulant matrices: A review. Technical report, Information Systems Laboratory, Department of Electrical Engineering, Stanford University, Stanford, California 94305, August 2002.
- [14] Simon Haykin. *Adaptive Filter Theory*. Prentice Hall, Upper Saddle River, NJ, 4th edition, 2002.
- [15] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 1985.
- [16] Mardia K. V., J.T. Kent and J.M. Bibby, *Multivariate Analysis*, Academic Press Ltd., 1979.
- [17] Mathpages on "Eigenvalue Problems and Matrix Invariants",

<http://www.mathpages.com/home/kmath128.htm>
- [18] Carl D. Meyer. Generalized inversion of modified matrices. *SIAM Journal of Applied Mathematics*, 24(3):315–323, May 1973.

- [19] Thomas P. Minka. Old and new matrix algebra useful for statistics, December 2000. Notes.
- [20] Daniele Mortari Ortho-Skew and Ortho-Sym Matrix Trigonometry *John Lee Junkins Astrodynamics Symposium*, AAS 03-265, May 2003. Texas A&M University, College Station, TX
- [21] L. Parra and C. Spence. Convolutional blind separation of non-stationary sources. In *IEEE Transactions Speech and Audio Processing*, pages 320–327, May 2000.
- [22] Kaare Brandt Petersen, Jiucang Hao, and Te-Won Lee. Generative and filtering approaches for overcomplete representations. *Neural Information Processing - Letters and Reviews*, vol. 8(1), 2005.
- [23] John G. Proakis and Dimitris G. Manolakis. *Digital Signal Processing*. Prentice-Hall, 1996.
- [24] Laurent Schwartz. *Cours d'Analyse*, volume II. Hermann, Paris, 1967. As referenced in [14].
- [25] Shayle R. Searle. *Matrix Algebra Useful for Statistics*. John Wiley and Sons, 1982.
- [26] G. Seber and A. Lee. *Linear Regression Analysis*. John Wiley and Sons, 2002.
- [27] S. M. Selby. *Standard Mathematical Tables*. CRC Press, 1974.
- [28] Inna Stainvas. Matrix algebra in differential calculus. Neural Computing Research Group, Information Engineering, Aston University, UK, August 2002. Notes.
- [29] P. P. Vaidyanathan. *Multirate Systems and Filter Banks*. Prentice Hall, 1993.
- [30] Max Welling. The Kalman Filter. Lecture Note.
- [31] Wikipedia on minors: "Minor (linear algebra)",

[http://en.wikipedia.org/wiki/Minor_\(linear_algebra\)](http://en.wikipedia.org/wiki/Minor_(linear_algebra))
- [32] Zhaoshui He, Shengli Xie, et al, "Convolutional blind source separation in frequency domain based on sparse representation", *IEEE Transactions on Audio, Speech and Language Processing*, vol.15(5):1551-1563, July 2007.
- [33] Karim T. Abou-Moustafa *On Derivatives of Eigenvalues and Eigenvectors of the Generalized Eigenvalue Problem*. McGill Technical Report, October 2010.
- [34] Mohammad Emtiyaz Khan *Updating Inverse of a Matrix When a Column is Added/Removed*. Emt CS,UBC February 27, 2008

Index

- Anti-symmetric, 54
- Block matrix, 46
- Chain rule, 15
- Cholesky-decomposition, 32
- Co-kurtosis, 34
- Co-skewness, 34
- Condition number, 62
- Cramers Rule, 29
- Derivative of a complex matrix, 24
- Derivative of a determinant, 8
- Derivative of a trace, 12
- Derivative of an inverse, 9
- Derivative of symmetric matrix, 15
- Derivatives of Toeplitz matrix, 16
- Dirichlet distribution, 37
- Eigenvalues, 30
- Eigenvectors, 30
- Exponential Matrix Function, 59
- Gaussian, conditional, 40
- Gaussian, entropy, 44
- Gaussian, linear combination, 41
- Gaussian, marginal, 40
- Gaussian, product of densities, 42
- Generalized inverse, 21
- Hadamard inequality, 52
- Hermitian, 48
- Idempotent, 49
- Kronecker product, 59
- LDL decomposition, 33
- LDM-decomposition, 33
- Linear regression, 28
- LU decomposition, 32
- Lyapunov Equation, 30
- Moore-Penrose inverse, 21
- Multinomial distribution, 37
- Nilpotent, 49
- Norm of a matrix, 61
- Norm of a vector, 61
- Normal-Inverse Gamma distribution, 37
- Normal-Inverse Wishart distribution, 39
- Orthogonal, 49
- Power series of matrices, 58
- Probability matrix, 55
- Pseudo-inverse, 21
- Schur complement, 41, 47
- Single entry matrix, 52
- Singular Valued Decomposition (SVD), 31
- Skew-Hermitian, 48
- Skew-symmetric, 54
- Stochastic matrix, 55
- Student-t, 37
- Sylvester's Inequality, 62
- Symmetric, 54
- Taylor expansion, 58
- Toeplitz matrix, 54
- Transition matrix, 55
- Trigonometric functions, 59
- Unipotent, 49
- Vandermonde matrix, 57
- Vec operator, 59, 60
- Wishart distribution, 38
- Woodbury identity, 18