

DATTORRO

CONVEX  
OPTIMIZATION



EUCLIDEAN  
DISTANCE  
GEOMETRY<sup>2ε</sup>

MeBoo

DATTORO

CONVEX

OPTIMIZATION

$\dagger$

EUCLIDEAN

DISTANCE

GEOMETRY

$2\epsilon$

MEBOO

# Convex Optimization

†

# Euclidean Distance Geometry $2\varepsilon$



*Mεβoo* Publishing

**Meboo Publishing USA**  
PO Box 12  
Palo Alto, California 94302

**Dattorro**, *Convex Optimization + Euclidean Distance Geometry*, second edition,  
*Mεβoo*, v2017.09.11.

ISBN 978-0-578-16140-2

This is version 2017.09.11: available in print at [Lulu.com](#), as conceived, in color.

cybersearch:

- I. semidefinite program
- II. rank constraint
- III. cardinality minimization
- IV. audio signal processing
- V. convex geometry
- VI. convex cones
- VII. distance matrix



with donations from **MATHWORKS, SIAM, and AMS**.

COVER ART DERIVES FROM NEW TOWN HALL, MUNICH



BY [SZE WAN](#)

This searchable electronic color pdfBook is click-navigable within the text by page, section, subsection, chapter, theorem, example, definition, cross reference, citation, equation, figure, table, and hyperlink. A pdfBook has no electronic copy protection and can be read and printed by most computers. The publisher hereby grants the right to reproduce this work in any format but limited to personal use.

for *Jennie Columba*



◇ *Antonio*



& *Sze Wan*

$$\texttt{EDM} = \mathbb{S}_h \cap \left( \mathbb{S}_c^\perp - \mathbb{S}_+ \right)$$

# Prelude

*The constant demands of my department and university and the ever increasing work needed to obtain funding have stolen much of my precious thinking time, and I sometimes yearn for the halcyon days of Bell Labs.*

—Steven Chu, Nobel laureate [89]

Convex Analysis is an emerging calculus of inequalities while Convex Optimization is its application. Analysis is inherently the domain of a mathematician while Optimization belongs to the engineer. A convex optimization problem is conventionally regarded as minimization of a convex objective function subject to an artificial convex domain imposed upon it by the problem constraints. The constraints comprise equalities and inequalities of convex functions whose simultaneous solution set generally constitutes the imposed convex domain: called *feasible set*.

It is easy to minimize a convex function over any convex subset of its domain because any local minimum must be a global minimum. But it is difficult to find the maximum of a convex function over some convex domain because there can be many local maxima; although this has practical application (Eternity II §4.7.0.0.15, §C.5), it is not a convex problem. Tremendous benefit accrues when a mathematical problem can be transformed to an equivalent convex optimization, primarily because any locally optimal solution is then guaranteed globally optimal.<sup>0.1</sup> An *optimal* solution is a best solution to the problem posed; a certificate can be obtained guaranteeing that no better solution exists.

---

<sup>0.1</sup>Solving a nonlinear system, for example, by instead solving an equivalent convex optimization problem is therefore highly preferable and what motivates *geometric programming*; a form of convex optimization invented in 1960s [64] [87] that has driven great advances in the electronic circuit design industry. [35, §4.7] [278] [445] [448] [113] [206] [215] [216] [217] [218] [219] [295] [296] [345]

Recognizing a problem as convex is an acquired skill; that being, to know when an objective function is convex and when constraints specify a convex feasible set. The challenge, which is indeed an art, is how to express difficult problems in a convex way: perhaps, problems previously believed nonconvex. Practitioners in the art of Convex Optimization engage themselves with discovery of which hard problems can be transformed into convex equivalents; because, once convex form of a problem is found, then a globally optimal solution is close at hand - the hard work is finished: Finding convex expression of a problem is itself, in a very real sense, its solution.

Yet, that skill acquired by understanding the geometry and application of Convex Optimization will remain more an art for some time to come; the reason being, there is generally no unique transformation of a given problem to its convex equivalent. This means, two researchers pondering the same problem are likely to formulate a convex equivalent differently; hence, one solution is likely different from the other although any convex combination of those two solutions remains optimal. Any presumption of only one right or correct solution becomes nebulous. Study of equivalence & sameness, uniqueness, and duality therefore pervade study of Optimization.

It can be difficult for the engineer to apply convex theory without an understanding of Analysis. These pages comprise my journal over a ten year period bridging gaps between engineer and mathematician; they constitute a translation, unification, and cohering of about four hundred papers, books, and reports from several different fields of mathematics and engineering. Although beacons of historical accomplishment are cited throughout, much of what is written here will not be found elsewhere. Care to detail, clarity, accuracy, consistency, and typography accompanies removal of ambiguity and verbosity out of respect for the reader. But the book is nonlinear in its presentation. Consequently there is much indexing, cross referencing, linkage to online sources, and background material provided in the text, footnotes, and appendices so as to be more self-contained and to provide understanding of fundamental concepts.

Looking toward the future, there remains much to be done in the area of machine computation if mathematical Optimization is to become fully embraced by the signal processing community. Wordlength of contemporary computers and numerical burdens upon them prohibit real time solution and accuracy sufficient to embed optimization problems within a recursive mathematical setting. When optimization problems constitute only intermediate solution to much larger problems, acquiring only a “few digits” accuracy can throw off subsequent dependent calculations. *Barrier* methods of solution are the principle obstacle to accuracy while *simplex* methods are the principle setback to speed. Novel, not hybrid, methods of solution are needed.

— Jon Dattorro  
Stanford, California  
2015

# Convex Optimization

## Euclidean Distance Geometry<sup>2ε</sup>

<b>1 Overview</b>	<b>19</b>
<b>2 Convex geometry</b>	<b>31</b>
2.1 Convex set . . . . .	31
2.2 Vectorized-matrix inner product . . . . .	42
2.3 Hulls . . . . .	50
2.4 Halfspace, Hyperplane . . . . .	58
2.5 Subspace representations . . . . .	69
2.6 Extreme, Exposed . . . . .	74
2.7 Cones . . . . .	77
2.8 Cone boundary . . . . .	85
2.9 Positive semidefinite (PSD) cone . . . . .	90
2.10 Conic independence (c.i.) . . . . .	111
2.11 When extreme means exposed . . . . .	115
2.12 Convex polyhedra . . . . .	116
2.13 Dual cone & generalized inequality . . . . .	122
<b>3 Geometry of convex functions</b>	<b>169</b>
3.1 Convex real and vector-valued function . . . . .	169
3.2 Practical norm functions, absolute value . . . . .	173
3.3 Powers, roots, and inverted functions . . . . .	180
3.4 Affine function . . . . .	182
3.5 Epigraph, Sublevel set . . . . .	185
3.6 Gradient . . . . .	192
3.7 First-order convexity condition, real function . . . . .	197
3.8 First-order convexity condition, vector-valued . . . . .	199
3.9 Second-order convexity condition, vector-valued . . . . .	200
3.10 Convex matrix-valued function . . . . .	201
3.11 First-order convexity condition, matrix-valued . . . . .	203
3.12 Epigraph of matrix-valued function, sublevel sets . . . . .	204
3.13 Second-order convexity condition, matrix-valued . . . . .	204
3.14 Quasiconvex . . . . .	206
3.15 Salient properties . . . . .	218
<b>4 Semidefinite programming</b>	<b>221</b>
4.1 Conic problem . . . . .	222
4.2 Framework . . . . .	228
4.3 Rank reduction . . . . .	238

4.4	Cardinality reduction . . . . .	244
4.5	Rank constraint by Convex Iteration . . . . .	247
4.6	Constraining cardinality . . . . .	271
4.7	Cardinality and rank constraint examples . . . . .	286
4.8	Quantum optimization . . . . .	324
4.9	Constraining rank of indefinite matrices . . . . .	329
4.10	Convex Iteration rank-1 . . . . .	333
4.11	Convex Iteration accelerant . . . . .	338
<b>5</b>	<b>Euclidean Distance Matrix</b>	<b>341</b>
5.1	EDM . . . . .	341
5.2	First metric properties . . . . .	342
5.3	$\exists$ fifth Euclidean metric property . . . . .	343
5.4	EDM definition . . . . .	346
5.5	Invariance . . . . .	370
5.6	Injectivity of $\mathbf{D}$ & unique reconstruction . . . . .	373
5.7	Embedding in affine hull . . . . .	378
5.8	Euclidean metric <i>versus</i> matrix criteria . . . . .	382
5.9	Bridge: Convex polyhedra to EDMs . . . . .	387
5.10	EDM-entry composition . . . . .	392
5.11	EDM indefiniteness . . . . .	395
5.12	List reconstruction . . . . .	399
5.13	Reconstruction examples . . . . .	403
5.14	Fifth property of Euclidean metric . . . . .	407
<b>6</b>	<b>Cone of distance matrices</b>	<b>415</b>
6.1	Defining EDM cone . . . . .	416
6.2	Polyhedral bounds . . . . .	418
6.3	$\sqrt{\text{EDM}}$ cone is not convex . . . . .	419
6.4	EDM definition in $\mathbf{1}\mathbf{1}^T$ . . . . .	419
6.5	Correspondence to PSD cone $\mathbb{S}_+^{N-1}$ . . . . .	426
6.6	Vectorization & projection interpretation . . . . .	430
6.7	A geometry of completion . . . . .	436
6.8	Dual EDM cone . . . . .	441
6.9	Theorem of the alternative . . . . .	453
6.10	Postscript . . . . .	453
<b>7</b>	<b>Proximity problems</b>	<b>455</b>
7.1	First prevalent problem: . . . . .	460
7.2	Second prevalent problem: . . . . .	468
7.3	Third prevalent problem: . . . . .	476
7.4	Conclusion . . . . .	483
<b>A</b>	<b>Linear algebra</b>	<b>485</b>
A.1	Main-diagonal $\delta$ operator, $\lambda$ , $\text{tr}$ , $\text{vec}$ , $\circ$ , $\otimes$ . . . . .	485
A.2	Semidefiniteness: domain of test . . . . .	488
A.3	Proper statements of positive semidefiniteness . . . . .	491
A.4	Schur complement . . . . .	499
A.5	Eigenvalue decomposition . . . . .	503
A.6	Singular value decomposition, SVD . . . . .	506
A.7	Zeros . . . . .	510

<b>B Simple matrices</b>	<b>515</b>
B.1 Rank-one matrix (dyad) . . . . .	515
B.2 Doublet . . . . .	519
B.3 Elementary matrix . . . . .	520
B.4 Auxiliary $V$ -matrices . . . . .	522
B.5 Orthomatrices . . . . .	525
B.6 Arrow matrix . . . . .	528
<b>C Some analytical optimal results</b>	<b>529</b>
C.1 Properties of infima . . . . .	529
C.2 Trace, singular and eigen values . . . . .	530
C.3 Orthogonal Procrustes problem . . . . .	535
C.4 Two-sided orthogonal Procrustes . . . . .	537
C.5 Quadratics . . . . .	540
<b>D Matrix calculus</b>	<b>543</b>
D.1 Gradient, Directional derivative, Taylor series . . . . .	543
D.2 Tables of gradients and derivatives . . . . .	558
<b>E Projection</b>	<b>567</b>
E.1 Idempotent matrices . . . . .	571
E.2 $I - P$ , Projection on algebraic complement . . . . .	574
E.3 Symmetric idempotent matrices . . . . .	575
E.4 Algebra of projection on affine subsets . . . . .	579
E.5 Projection examples . . . . .	580
E.6 Vectorization interpretation . . . . .	587
E.7 Projection on matrix subspaces . . . . .	592
E.8 Range Rowspace interpretation . . . . .	594
E.9 Projection on convex set . . . . .	595
E.10 Alternating projection . . . . .	606
<b>F Notation, Definitions, Glossary</b>	<b>621</b>
<b>Bibliography</b>	<b>637</b>
<b>Index</b>	<b>659</b>

# List of Tables

<b>2 Convex geometry</b>	
Table 2.9.2.3.1, rank <i>versus</i> dimension of $\mathbb{S}_+^3$ faces	97
Table 2.10.0.0.1, maximum number of c.i. directions	111
Cone Table 1	151
Cone Table S	152
Cone Table A	153
Cone Table 1*	157
<b>4 Semidefinite programming</b>	
Faces of $\mathbb{S}_+^3$ corresponding to faces of $\mathcal{S}_+^3$	226
Quantum impulse	326
Quantum step	328
Quantum AND function	328
<b>5 Euclidean Distance Matrix</b>	
Précis 5.7.2: affine dimension in terms of rank	381
<b>B Simple matrices</b>	
Auxiliary V-matrix Table B.4.4	524
<b>D Matrix calculus</b>	
Table D.2.1, algebraic gradients and derivatives	559
Table D.2.2, trace Kronecker gradients	560
Table D.2.3, trace gradients and derivatives	561
Table D.2.4, logarithmic determinant gradients, derivatives	563
Table D.2.5, determinant gradients and derivatives	564
Table D.2.6, logarithmic derivatives	564
Table D.2.7, exponential gradients and derivatives	565

# List of Figures

<b>1</b>	<b>Overview</b>	<b>19</b>
1	Sigma delta quantizer . . . . .	20
2	Room geometry estimation by first acoustic reflections . . . . .	20
3	<i>Orion nebula</i> . . . . .	21
4	Application of trilateration is localization . . . . .	22
5	Molecular conformation . . . . .	23
6	Facial recognition . . . . .	24
7	<i>Swiss roll</i> . . . . .	25
8	USA map reconstruction . . . . .	26
9	Honeycomb, Hexabenzocoronene molecule . . . . .	27
10	Robotic vehicles . . . . .	28
11	Reconstruction of David . . . . .	29
12	David by distance geometry . . . . .	29
<b>2</b>	<b>Convex geometry</b>	<b>31</b>
13	Slab . . . . .	33
14	Open, closed, convex sets . . . . .	35
15	Intersection of line with boundary . . . . .	36
16	Tangentials . . . . .	38
17	Inverse image . . . . .	41
18	Inverse image under linear map . . . . .	41
19	<i>Tesseract</i> . . . . .	44
20	Linear injective mapping of Euclidean body . . . . .	45
21	Linear noninjective mapping of Euclidean body . . . . .	46
22	Convex hull of a random list of points . . . . .	50
23	Hulls . . . . .	52
24	Two Fantopes . . . . .	54
25	Nuclear Norm Ball . . . . .	55
26	Convex hull of rank-1 matrices . . . . .	56
27	A simplicial cone . . . . .	59
28	Hyperplane illustrated $\partial\mathcal{H}$ is a partially bounding line . . . . .	60
29	Hyperplanes in $\mathbb{R}^2$ . . . . .	62
30	Affine independence . . . . .	64
31	$\{z \in \mathcal{C} \mid a^T z = \kappa_i\}$ . . . . .	65
32	Hyperplane supporting closed set . . . . .	66
33	Minimizing hyperplane over affine subset in nonnegative orthant . . . . .	72
34	Maximizing hyperplane over convex set . . . . .	73
35	Closed convex set illustrating exposed and extreme points . . . . .	78

36	Two-dimensional nonconvex cone . . . . .	78
37	Nonconvex cone made from lines . . . . .	79
38	Nonconvex cone is convex cone boundary . . . . .	79
39	Union of convex cones is nonconvex cone . . . . .	79
40	Truncated nonconvex cone $\mathcal{X}$ . . . . .	80
41	Cone exterior is convex cone . . . . .	80
42	Not a cone . . . . .	81
43	Minimum element, Minimal element . . . . .	83
44	$\mathcal{K}$ is a pointed polyhedral cone not full-dimensional . . . . .	86
45	Exposed and extreme directions . . . . .	89
46	Positive semidefinite cone . . . . .	92
47	Convex Schur-form set . . . . .	93
48	Projection of truncated PSD cone . . . . .	95
49	Circular cone showing axis of revolution . . . . .	103
50	Circular section . . . . .	104
51	Polyhedral inscription . . . . .	106
52	Conically (in)dependent vectors . . . . .	112
53	Pointed six-faceted polyhedral cone and its dual . . . . .	113
54	Minimal set of generators for halfspace about origin . . . . .	115
55	Venn diagram for cones and polyhedra . . . . .	117
56	Range form polyhedron . . . . .	118
57	Simplex . . . . .	120
58	Two views of a simplicial cone and its dual . . . . .	121
59	Two equivalent constructions of dual cone . . . . .	123
60	Dual polyhedral cone construction by right angle . . . . .	124
61	Orthogonal cones . . . . .	126
62	Blades $\mathcal{K}$ and $\mathcal{K}^*$ . . . . .	127
63	$\mathcal{K}$ is a halfspace about the origin . . . . .	128
64	Iconic primal and dual objective functions . . . . .	129
65	Membership w.r.t $\mathcal{K}$ and orthant . . . . .	137
66	Shrouded polyhedral cone . . . . .	142
67	Simplicial cone $\mathcal{K}$ in $\mathbb{R}^2$ and its dual . . . . .	146
68	Monotone nonnegative cone $\mathcal{K}_{\mathcal{M}+}$ and its dual . . . . .	154
69	Monotone cone $\mathcal{K}_{\mathcal{M}}$ and its dual . . . . .	155
70	Two views of monotone cone $\mathcal{K}_{\mathcal{M}}$ and its dual . . . . .	156
71	First-order optimality condition . . . . .	159
<b>3</b>	<b>Geometry of convex functions</b> . . . . .	<b>169</b>
72	Convex functions having unique minimizer . . . . .	170
73	Minimum/Minimal element, dual cone characterization . . . . .	172
74	1-norm ball $\mathcal{B}_1$ from compressed sensing/compressive sampling . . . . .	175
75	Cardinality minimization, signed <i>versus</i> unsigned variable . . . . .	176
76	1-norm variants . . . . .	176
77	Affine function . . . . .	184
78	Supremum of affine functions . . . . .	185
79	Epigraph . . . . .	185
80	Log function constraint . . . . .	191
81	Quadratic bowl and 1-norm gradients in $\mathbb{R}^2$ evaluated on grid . . . . .	192
82	Quadratic function convexity in terms of its gradient . . . . .	198
83	Contour plot of convex real function at selected levels . . . . .	199
84	Taxicab distance on nonuniform rectangular grid . . . . .	203

85	Iconic quasiconvex function . . . . .	207
86	Quasiconcave monotonic function $xu$ . . . . .	208
87	Operational Amplifier implementation of third-order filter having a zero . . . . .	210
88	Mason flowgraph for operational amplifier arbitrary magnitude filter circuit . . . . .	211
89	Bisection method linearity . . . . .	214
90	Arbitrary magnitude analog filter design . . . . .	215
91	Sum of convex functions . . . . .	218
<b>4</b>	<b>Semidefinite programming</b>	<b>221</b>
92	Venn diagram of convex program classes . . . . .	224
93	Visualizing positive semidefinite cone in high dimension . . . . .	225
94	Primal/Dual transformations . . . . .	232
95	Projection <i>versus</i> convex iteration . . . . .	249
96	Trace heuristic . . . . .	250
97	Sensor-network localization . . . . .	253
98	2-lattice of sensors and anchors for localization example . . . . .	255
99	3-lattice of sensors and anchors for localization example . . . . .	256
100	4-lattice of sensors and anchors for localization example . . . . .	257
101	5-lattice of sensors and anchors for localization example . . . . .	258
102	Uncertainty ellipsoids orientation and eccentricity . . . . .	259
103	2-lattice localization solution . . . . .	261
104	3-lattice localization solution . . . . .	262
105	4-lattice localization solution . . . . .	262
106	5-lattice localization solution . . . . .	263
107	10-lattice localization solution . . . . .	263
108	100 randomized noiseless sensor localization . . . . .	264
109	100 randomized sensors localization . . . . .	265
110	Nonnegative spectral factorization . . . . .	268
111	Regularization curve for convex iteration . . . . .	270
112	1-norm heuristic . . . . .	272
113	Sparse sampling theorem . . . . .	275
114	Signal dropout . . . . .	278
115	Signal dropout reconstruction . . . . .	279
116	Simplex with intersecting line problem in compressed sensing . . . . .	281
117	Geometric interpretations of sparse-sampling constraints . . . . .	283
118	Permutation matrix column-norm and column-sum constraint . . . . .	289
119	MAX CUT problem . . . . .	295
120	Shepp-Logan phantom . . . . .	299
121	MRI radial sampling pattern in Fourier domain . . . . .	302
122	Aliased phantom . . . . .	303
123	Neighboring-pixel stencil on Cartesian grid . . . . .	305
124	Differentiable almost everywhere . . . . .	306
125	<i>Eternity II</i> . . . . .	308
126	<i>Eternity II</i> game-board grid . . . . .	310
127	<i>Eternity II</i> demo-game piece illustrating edge-color ordering . . . . .	311
128	<i>Eternity II</i> vectorized demo-game-board piece descriptions . . . . .	312
129	<i>Eternity II</i> difference $\Delta$ and boundary coefficient $\beta$ construction . . . . .	313
130	<i>Eternity II</i> composite variable matrix sparsity pattern . . . . .	315
131	<i>Eternity II</i> problem visualization in three dimensions . . . . .	320
132	<i>Eternity II</i> permutation polyhedron vertices visualization on sphere . . . . .	321
133	<i>Chimera</i> topology for D:Wave 1152-qubit chip . . . . .	323

134	D:Wave <i>Chimera</i> chip layout	325
135	MIT logo	331
136	One-pixel camera	331
137	One-pixel camera - compression estimates	332
138	Convergence of Singular Value Decomposition by Convex Iteration	337
139	Straight line through three direction vectors by midpoint fit	338
<b>5</b>	<b>Euclidean Distance Matrix</b>	<b>341</b>
140	Convex hull of three points	342
141	Complete dimensionless <i>EDM graph</i>	344
142	Fifth Euclidean metric property	345
143	<i>Fermat point</i>	352
144	Arbitrary hexagon in $\mathbb{R}^3$	353
145	Kissing number	354
146	<i>Trilateration</i>	358
147	This EDM graph provides unique isometric reconstruction	361
148	Two sensors $\bullet$ and three anchors $\circ$	361
149	Two discrete linear trajectories of sensors	362
150	Coverage in cellular telephone network	365
151	Contours of equal signal power	365
152	Depiction of molecular conformation	366
153	Square diamond	372
154	Orthogonal complements in $\mathbb{S}^N$ abstractly oriented	374
155	Elliptope $\mathcal{E}^3$	388
156	Elliptope $\mathcal{E}^2$ interior to $\mathbb{S}_+^2$	389
157	Smallest eigenvalue of $-V_{\mathcal{N}}^T D V_{\mathcal{N}}$	393
158	Some entrywise EDM compositions	393
159	Map of United States of America	402
160	Largest ten eigenvalues of $-V_{\mathcal{N}}^T O V_{\mathcal{N}}$	404
161	<i>Relative-angle inequality tetrahedron</i>	409
162	Nonsimplicial pyramid in $\mathbb{R}^3$	412
<b>6</b>	<b>Cone of distance matrices</b>	<b>415</b>
163	Relative boundary of cone of Euclidean distance matrices	417
164	Example of $V_{\mathcal{X}}$ selection to make an EDM	421
165	Vector $V_{\mathcal{X}}$ spirals	423
166	Three views of translated negated elliptope	429
167	Halfline $\mathcal{T}$ on PSD cone boundary	432
168	Vectorization and projection interpretation example	433
169	Intersection of EDM cone with hyperplane	435
170	Neighborhood graph	437
171	<i>Trefoil knot</i> untied	438
172	<i>Trefoil ribbon</i>	440
173	Orthogonal complement of geometric center subspace	444
174	EDM cone construction by flipping PSD cone	445
175	Decomposing member of polar EDM cone	448
176	Ordinary dual EDM cone projected on $\mathbb{S}_h^3$	452

<b>7 Proximity problems</b>	<b>455</b>
177 Pseudo-Venn diagram for EDM . . . . .	457
178 Elbow placed in path of projection . . . . .	457
179 Convex envelope . . . . .	471
<b>A Linear algebra</b>	<b>485</b>
180 Geometrical interpretation of full SVD . . . . .	508
<b>B Simple matrices</b>	<b>515</b>
181 Four fundamental subspaces for any dyad . . . . .	516
182 Four fundamental subspaces for doublet . . . . .	519
183 Four fundamental subspaces for elementary matrix . . . . .	520
184 Gimbal . . . . .	526
185 Arrow matrix . . . . .	528
<b>D Matrix calculus</b>	<b>543</b>
186 Convex quadratic bowl in $\mathbb{R}^2 \times \mathbb{R}$ . . . . .	551
<b>E Projection</b>	<b>567</b>
187 Action of pseudoinverse . . . . .	568
188 Nonorthogonal projection of $x \in \mathbb{R}^3$ on $\mathcal{R}(U) = \mathbb{R}^2$ . . . . .	573
189 Biorthogonal expansion of point $x \in \text{aff } \mathcal{K}$ . . . . .	581
190 Linear regression <i>versus</i> principal component analysis . . . . .	585
191 Dual interpretation of projection on convex set . . . . .	597
192 Projection on orthogonal complement . . . . .	599
193 Projection on dual cone . . . . .	601
194 Projection product on convex set in subspace . . . . .	605
195 von Neumann-style projection of point $b$ . . . . .	607
196 Alternating projection on two halfspaces . . . . .	608
197 Distance, feasibility, optimization . . . . .	609
198 Alternating projection on nonnegative orthant and hyperplane . . . . .	611
199 Geometric convergence of iterates in norm . . . . .	612
200 Distance between PSD cone and iterate in $\mathcal{A}$ . . . . .	615
201 Dykstra's alternating projection algorithm . . . . .	616
202 Polyhedral normal cones . . . . .	617
203 Normal cone to ellotope . . . . .	618
204 Normal-cone progression . . . . .	620



# Chapter 1

## Overview

### Convex Optimization Euclidean Distance Geometry $2\varepsilon$

*People are so afraid of convex analysis.*

—Claude Lemaréchal, 2003

In layman's terms, the mathematical science of Optimization is a study of how to make good choices when confronted with conflicting requirements and demands. Optimization is a relatively new wisdom, historically, that can represent balance of real things. The qualifier *convex* means: when an optimal solution is found, then it is guaranteed to be a best solution; there is no better choice.

Any convex optimization problem has geometric interpretation. If a given optimization problem can be transformed to a convex equivalent, then this interpretive benefit is acquired. That is a powerful attraction: the ability to visualize geometry of an optimization problem. Conversely, recent advances in geometry and in graph theory hold convex optimization within their proofs' core. [457] [356]

This book is about convex optimization, convex geometry (with particular attention to distance geometry), and nonconvex, combinatorial, and geometrical problems that can be relaxed or transformed into convexity. A virtual flood of new applications follows by epiphany that many problems, presumed nonconvex, can be so transformed: [11] [12] [35, §4.3, p.316-322] [63] [102] [169] [172] [309] [334] [342] [402] [403] [453] [457] e.g, sigma delta analog-to-digital (A/D) audio converter antialiasing (Figure 1).

Euclidean distance geometry is, fundamentally, a determination of point conformation (configuration, relative position or location) by inference from interpoint distance information. By *inference* we mean: e.g, given only distance information, determine whether there corresponds a *realizable* conformation of points; a *list* of points in some dimension that attains the given interpoint distances. Each point may represent simply location or, abstractly, any entity expressible as a vector in finite-dimensional Euclidean space; e.g, distance geometry of music [120].

It is a common misconception to presume that some desired point conformation cannot be recovered in absence of complete interpoint distance information. We might, for example, want to realize a constellation given only interstellar distance (or, equivalently, parsecs from our Sun and relative angular measurement; the Sun as vertex to two distant stars); called *stellar cartography*, an application evoked by Figure 3. At first it may seem

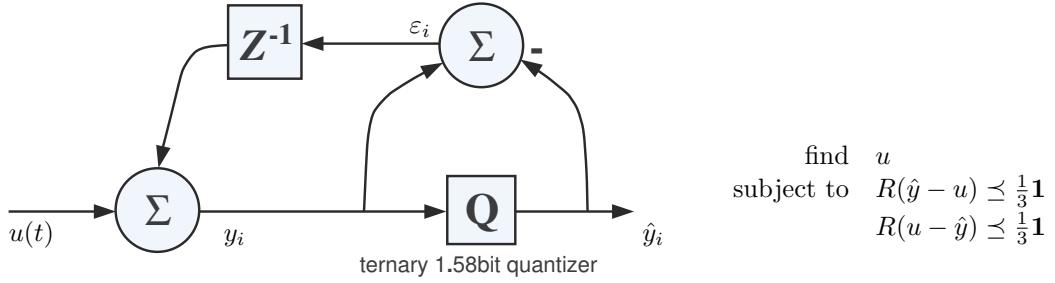


Figure 1: Multibit sigma delta quantization is predominant technology for analog to digital audio signal conversion. [2, p.6] Input signal  $u(t)$  is continuous. Delay  $z^{-1}$  here is analog, perhaps implemented by sample/hold circuit at MHz rate of  $\hat{y}_i$  samples. Observing vector  $\hat{y}$ , signal  $u$  can be reconstructed by finding a point feasible to the set of linear inequalities representing this coarse quantizer recursion.  $R$  is a lower triangular matrix of ones. [110]

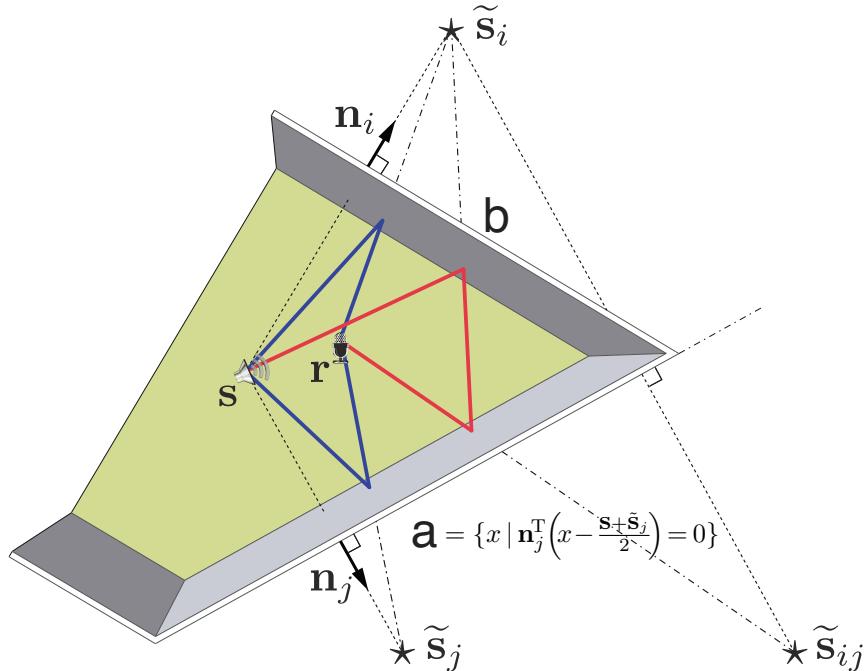


Figure 2: [132] [322] [129] Dokmanić & Parhizkar *et alii* discover an audio signal processing application of Euclidean distance matrices to room geometry estimation by discerning first acoustic reflections of stationary sound source  $s$ . Locations of source and phantom  $\star$  sources  $\tilde{s}_i$  and  $\tilde{s}_j$  are ascertained by measuring arrival times of first echoes (blue) at multiple microphone receivers. (Only one receiver  $r$  is illustrated. Second reflection (red) phantom  $\tilde{s}_{ij}$  ignored.) Phantom location is invariant to receiver position. All interpoint distances among receivers are known. Once source and phantoms are localized, normals  $\mathbf{n}_j$  and  $\mathbf{n}_i$  respectively identify truncated hyperplanes (walls)  $\mathbf{a}$  and  $\mathbf{b}$  bisecting perpendicular line segment connecting source  $s$  to a phantom.



Figure 3: *Orion nebula*. (Astrophotography by [Massimo Robberto](#).)

that  $O(N^2)$  data is required, yet there are many circumstances where this can be reduced to  $O(N)$ .

If we agree that a set of points may have a shape (three points can form a triangle and its interior, for example, four points a tetrahedron), then we can ascribe *shape* of a set of points to their convex hull. It should be apparent: from distance, these shapes can be determined only to within a *rigid transformation* (rotation, reflection, translation).

Absolute position information is generally lost, given only distance information, but we can determine the smallest possible dimension in which an unknown list of points can exist; that attribute is their *affine dimension* (a triangle in any ambient space has affine dimension 2, for example). In circumstances where stationary reference points are also provided, it becomes possible to determine absolute position or location; *e.g.*, Figure 4.

Geometric problems involving distance between points can sometimes be reduced to convex optimization problems. Mathematics of this combined study of geometry and optimization is rich and deep. Its application has already proven invaluable discerning organic *molecular conformation* by measuring interatomic distance along covalent bonds; *e.g.*, Figure 5. [96] [392] [157] [49] Many disciplines have already benefitted and simplified consequent to this theory; *e.g.*, distance based *pattern recognition* (Figure 6), *localization* in wireless sensor networks [50] [451] [48] by measurement of intersensor distance along channels of communication, *wireless location* of a radio-signal source such as cell phone by multiple measurements of signal strength, the *global positioning system* (GPS), *multidimensional scaling* (§5.12) which is a numerical representation of qualitative data by finding a low-dimensional scale, and audio signal processing: ultrasound tomography, room geometry estimation (Figure 2), and perhaps dereverberation by localization of phantom sound sources [130] [129] [132]. [131]

Euclidean distance geometry provides some foundation for *artificial intelligence*. Together with convex optimization, distance geometry has found application to:

- *machine learning* by discerning naturally occurring manifolds in:
  - Euclidean bodies (Figure 7, §6.7.0.0.1)
  - Fourier spectra of kindred utterances [239]
  - photographic image sequences [435]

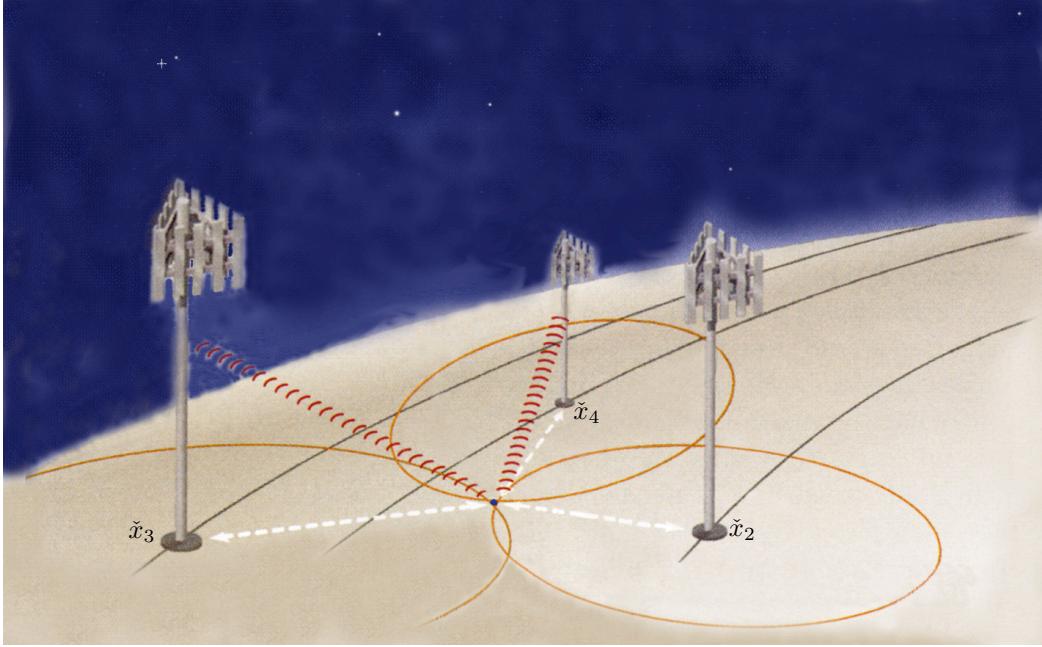


Figure 4: Application of trilateration (§5.4.2.2.8) is localization (determining position) of a radio signal source in 2 dimensions; more commonly known by radio engineers as the process “triangulation”. In this scenario, anchors  $\check{x}_2, \check{x}_3, \check{x}_4$  are illustrated as fixed antennae. [235] The radio signal source (a sensor  $\bullet x_1$ ) anywhere in affine hull of three antenna bases can be uniquely localized by measuring distance to each (dashed white arrowed line segments). Ambiguity of lone distance measurement to sensor is represented by circle about each antenna. Trilateration is expressible as a semidefinite program; hence, a convex optimization problem. [357]

- *robotics*; e.g., automated manufacturing, and autonomous navigation of vehicles maneuvering in formation (Figure 10).

## by chapter

We study the many manifestations and representations of pervasive convex Euclidean bodies. In particular, we make convex polyhedra, cones, and dual cones visceral through illustration in **chapter 2 Convex geometry** where geometric relation of polyhedral cones to nonorthogonal bases (biorthogonal expansion) is examined. It is shown that coordinates are unique in any conic system whose basis cardinality equals or exceeds spatial dimension; for high cardinality, a new definition of *conic coordinate* is provided in Theorem 2.13.13.0.1. The conic analogue to linear independence, called *conic independence*, is introduced as a tool for study, analysis, and manipulation of cones; a natural extension and next logical step in progression: linear, affine, conic. We explain conversion between halfspace- and vertex-description of a convex cone, we motivate the dual cone and provide formulae for finding it, and we show how first-order optimality conditions or alternative systems of linear inequality or *linear matrix inequality* can be explained by *dual generalized inequalities* with respect to convex cones. Arcane theorems of alternative generalized inequality are, in fact, simply derived from cone *membership relations*; generalizations of algebraic *Farkas’ lemma* translated to geometry of convex cones.

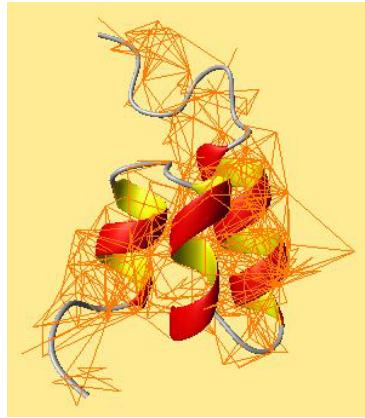


Figure 5: [214] [134] Distance data collected via nuclear magnetic resonance (NMR) helped render this three-dimensional depiction of a [protein molecule](#). At the beginning of the 1980s, Kurt Wüthrich [Nobel laureate] developed an idea about how NMR could be extended to cover biological molecules such as proteins. He invented a systematic method of pairing each NMR signal with the right hydrogen nucleus (proton) in the macromolecule. The method is called sequential assignment and is today a cornerstone of all NMR structural investigations. He also showed how it was subsequently possible to determine pairwise distances between a large number of hydrogen nuclei and use this information with a mathematical method based on distance-geometry to calculate a three-dimensional structure for the molecule. [440] [209] –[313]

Any convex optimization problem can be visualized geometrically. Desire to visualize in high dimension [[Sagan, \*Cosmos – The Edge of Forever\*, 22:55'](#)] is deeply embedded in the [mathematical psyche](#). [1] Chapter 2 provides tools to make visualization easier, and we teach how to visualize in high dimension. The concepts of face, extreme point, and extreme direction of a convex Euclidean body are explained here; crucial to understanding convex optimization. How to find the smallest face of any closed convex cone, containing convex set  $\mathcal{C}$ , is divulged; later shown to have practical application to presolving convex programs. The convex cone of positive semidefinite matrices, in particular, is studied in depth:

- We interpret, for example, inverse image of the positive semidefinite cone under affine transformation. (Example 2.9.1.0.2)
- Subsets of the positive semidefinite cone, discriminated by rank exceeding some lower bound, are convex. In other words, high-rank subsets of the positive semidefinite cone boundary united with its interior are convex. (Theorem 2.9.2.9.3) There is a closed form for projection on those convex subsets.
- The positive semidefinite cone is a circular cone in low dimension; *Geršgorin discs* specify inscription of a polyhedral cone into it. (Figure 51)

**Chapter 3 Geometry of convex functions** observes Fenchel's analogy between convex sets and functions: We explain, for example, how the real affine function relates to convex functions as the hyperplane relates to convex sets. Partly a toolbox of practical useful convex functions and a cookbook for optimization problems, methods are drawn from the appendices about matrix calculus for determining convexity and discerning geometry.

**Chapter 4. Semidefinite programming** has recently emerged to prominence because it admits a new class of problem previously unsolvable by convex optimization

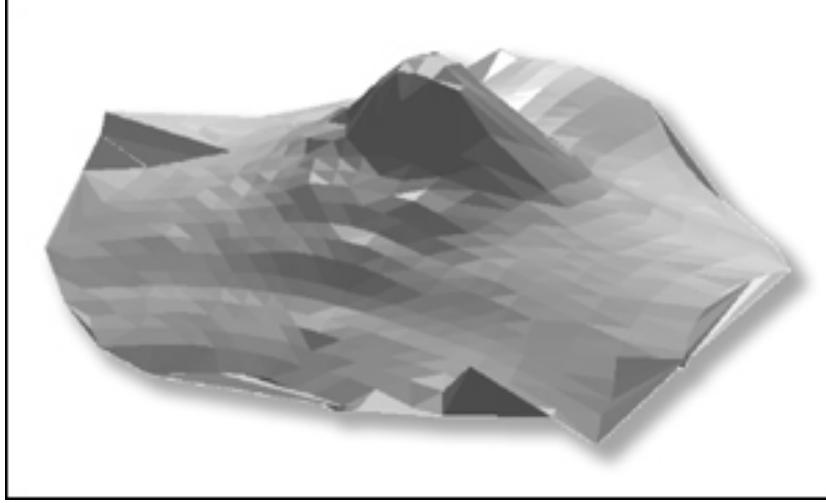


Figure 6: This coarsely discretized triangulated algorithmically flattened human face (made by Kimmel & the Bronsteins [253]) represents a stage in machine recognition of human identity; called *facial recognition*. Distance geometry is applied to determine discriminating-features.

*techniques, and because it theoretically subsumes other convex techniques: linear programming and quadratic programming and second-order cone programming.* –p.223 Semidefinite programming is reviewed with particular attention to optimality conditions for prototypical primal and dual problems, their interplay, and a perturbation method for rank reduction of optimal solutions (extant but not well known). *Positive definite Farkas' lemma* is derived, and we also show how to determine if a feasible set belongs exclusively to a positive semidefinite cone boundary. An arguably good three-dimensional polyhedral analogue to the positive semidefinite cone of  $3 \times 3$  symmetric matrices is introduced: a new tool for visualizing coexistence of low- and high-rank optimal solutions in six isomorphic dimensions and a mnemonic aid for understanding semidefinite programs. We find a minimal cardinality Boolean solution to an instance of  $Ax = b$ :

$$\begin{aligned} & \underset{x}{\text{minimize}} && \|x\|_0 \\ & \text{subject to} && Ax = b \\ & && x_i \in \{0, 1\}, \quad i=1 \dots n \end{aligned} \tag{726}$$

The *sensor-network localization* problem is solved in any dimension in this chapter. We introduce a method of *convex iteration* for constraining rank in the form  $\text{rank } G \leq \rho$  and cardinality in the form  $\text{card } x \leq k$ . Cardinality minimization is applied to a discrete image-gradient of the Shepp-Logan phantom, from Magnetic Resonance Imaging (MRI) in the field of medical imaging, for which we find a new lower bound of 1.9% cardinality. We show how to handle polynomial constraints, and how to transform a rank-constrained problem to a rank-1 problem.

The EDM is studied in **chapter 5 Euclidean Distance Matrix**; its properties and relationship to both positive semidefinite and Gram matrices. We relate the EDM to the four classical properties of Euclidean metric; thereby, observing existence of an infinity of properties of the Euclidean metric beyond triangle inequality. We proceed by deriving the fifth Euclidean metric property and then explain why furthering this endeavor is inefficient because the ensuing criteria (while describing polyhedra in angle or area, volume, content, and so on *ad infinitum*) grow linearly in complexity and number with problem size.

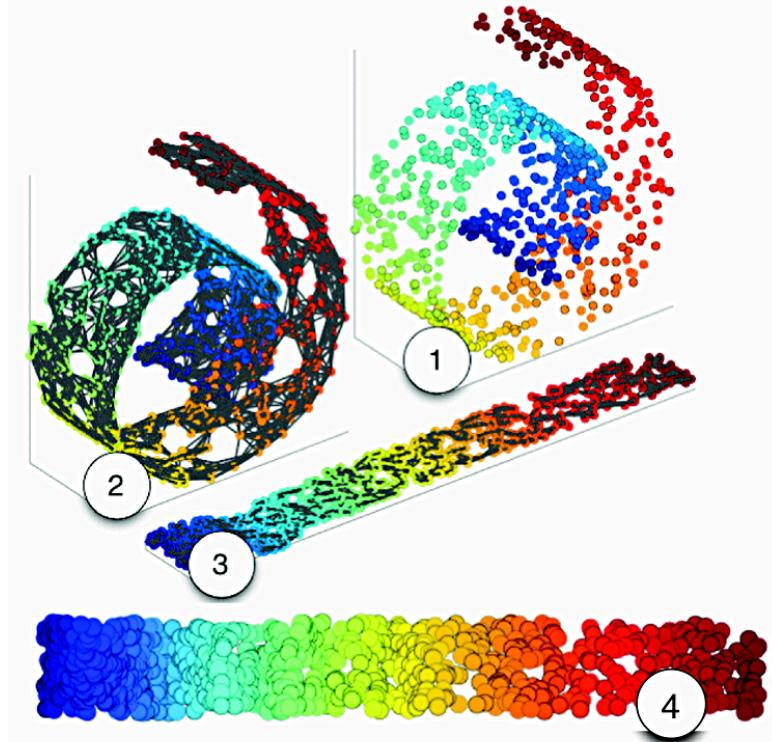


Figure 7: *Swiss roll*, Weinberger & Saul [435]. The problem of manifold learning, illustrated for  $N = 800$  data points sampled from a “Swiss roll” ①. A discretized manifold is revealed by connecting each data point and its  $k=6$  nearest neighbors ②. An unsupervised learning algorithm unfolds the Swiss roll while preserving the local geometry of nearby data points ③. Finally, the data points are projected onto the two-dimensional subspace that maximizes their variance, yielding a faithful embedding of the original manifold ④.

Reconstruction methods are explained and applied to a map of the United States; e.g., Figure 8. We also experimentally test a conjecture of Borg & Groenen by reconstructing a distorted but recognizable isotonic map of the USA using only ordinal (comparative) distance data: Figure 159e-f. We demonstrate an elegant method for including dihedral (or *torsion*) angle constraints into a molecular conformation problem. We explain why *trilateration* (a.k.a *localization*) is a convex optimization problem. We show how to recover relative position given incomplete interpoint distance information, and how to pose EDM problems or transform geometrical problems to convex optimizations; e.g., *kissing number* of packed spheres about a central sphere (solved in  $\mathbb{R}^3$  by Isaac Newton).

The set of all Euclidean distance matrices forms a pointed closed convex cone called the *EDM cone*:  $\mathbb{EDM}^N$ . We offer a new proof of Schoenberg’s seminal characterization of EDMs:

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} -V_{\mathcal{N}}^T D V_{\mathcal{N}} \succeq 0 \\ D \in \mathbb{S}_h^N \end{cases} \quad (1052)$$

Our proof relies on fundamental geometry; assuming, any EDM must correspond to a list of points contained in some polyhedron (possibly at its vertices) and *vice versa*. It is known, but not obvious, this *Schoenberg criterion* implies nonnegativity of the EDM entries; proved herein.

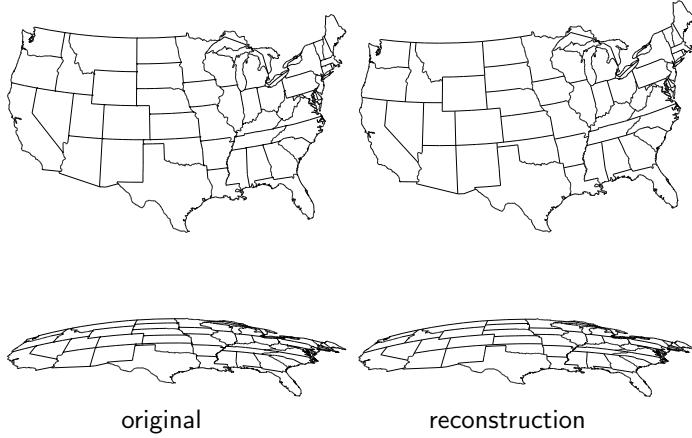


Figure 8: (confer Figure 159) About five thousand points along borders constituting United States were used to create an exhaustive matrix of interpoint distance for each and every pair of points in an ordered set (a *list*); called *Euclidean distance matrix*. From that noiseless distance information, it is easy to reconstruct this nonconvex map exactly via Schoenberg criterion (1052). (§5.13.1.0.1) Map reconstruction is exact (to within a rigid transformation) given any number of interpoint distances; the greater the number of distances, the greater the detail (as it is for all conventional map preparation).

We characterize eigenvalue spectrum of an EDM, then devise a polyhedral spectral cone for determining membership of a given matrix (in Cayley-Menger form) to the convex cone of Euclidean distance matrices; *id est*, a matrix is an EDM if and only if its nonincreasingly ordered vector of eigenvalues belongs to a polyhedral spectral cone for  $\mathbb{EDM}^N$

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} \lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix}\right) \in \left[\begin{array}{c} \mathbb{R}_+^N \\ \mathbb{R}_- \end{array}\right] \cap \partial\mathcal{H} \\ D \in \mathbb{S}_h^N \end{cases} \quad (1270)$$

We will see: spectral cones are not unique.

In **chapter 6 Cone of distance matrices** we explain a geometric relationship between the cone of Euclidean distance matrices, two positive semidefinite cones, and the ellipope. We illustrate geometric requirements, in particular, for projection of a given matrix on a positive semidefinite cone that establish its membership to the EDM cone. The faces of the EDM cone are described, but still open is the question whether all its faces are exposed as they are for the positive semidefinite cone.

The *Schoenberg criterion*,

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} -V_N^T D V_N \in \mathbb{S}_+^{N-1} \\ D \in \mathbb{S}_h^N \end{cases} \quad (1052)$$

for identifying a Euclidean distance matrix, is revealed to be a discretized *membership relation* (*dual generalized inequalities*, a new Farkas'-like lemma) between the EDM cone and its ordinary dual:  $\mathbb{EDM}^{N^*}$ . A matrix criterion for membership to the dual EDM cone is derived that is simpler than the Schoenberg criterion:

$$D^* \in \mathbb{EDM}^{N^*} \Leftrightarrow \delta(D^* \mathbf{1}) - D^* \succeq 0 \quad (1420)$$

There is a concise equality, relating the convex cone of Euclidean distance matrices to the positive semidefinite cone, apparently overlooked in the literature; an equality between

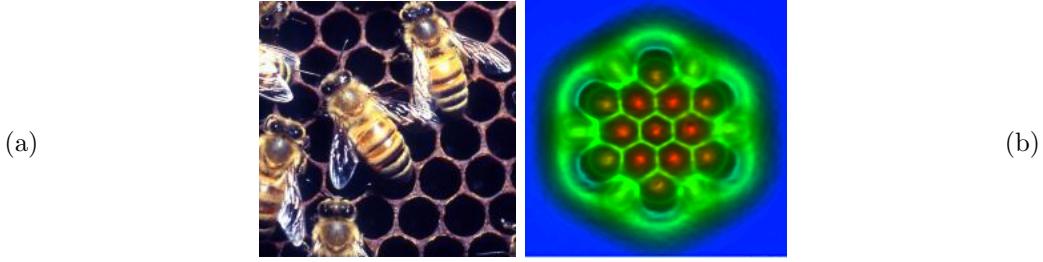


Figure 9: (a) These bees construct a honeycomb by solving a convex optimization problem (§5.4.2.2.6). The most dense packing of identical spheres about a central sphere in 2 dimensions is 6. Sphere centers describe a regular lattice. (b) A hexabenzocoronene molecule (diameter: 1.4nm) imaged by noncontact atomic force microscopy using a microscope tip terminated with a single carbon monoxide molecule. The carbon-carbon bonds in the imaged molecule appear with different contrast and apparent lengths. Based on these disparities, the bond orders and lengths of the individual bonds can be distinguished. (Image by Leo Gross.)

two large convex Euclidean bodies:

$$\text{EDM}^N = \mathbb{S}_h^N \cap \left( \mathbb{S}_c^{N\perp} - \mathbb{S}_+^N \right) \quad (1414)$$

Seemingly innocuous problems in terms of point position  $x_i \in \mathbb{R}^n$  like

$$\underset{\{x_i\}}{\text{minimize}} \sum_{i,j \in \mathcal{I}} (\|x_i - x_j\| - h_{ij})^2 \quad (1454)$$

$$\underset{\{x_i\}}{\text{minimize}} \sum_{i,j \in \mathcal{I}} (\|x_i - x_j\|^2 - h_{ij})^2 \quad (1455)$$

are difficult to solve. So, in **chapter 7 Proximity problems**, we instead explore methods of their solution by transformation to a few fundamental and prevalent Euclidean distance matrix proximity problems; the problem of finding that distance matrix closest, in some sense, to a given matrix  $H = [h_{ij}]$ :

$$\begin{array}{ll} \underset{D}{\text{minimize}} & \| -V(D - H)V \|^2_F \\ \text{subject to} & \text{rank } VDV \leq \rho \\ & D \in \text{EDM}^N \end{array} \quad \begin{array}{ll} \underset{\sqrt[D]{D}}{\text{minimize}} & \| \sqrt[D]{D} - H \|^2_F \\ \text{subject to} & \text{rank } VDV \leq \rho \\ & \sqrt[D]{D} \in \sqrt{\text{EDM}^N} \end{array} \quad (1456)$$

$$\begin{array}{ll} \underset{D}{\text{minimize}} & \| D - H \|^2_F \\ \text{subject to} & \text{rank } VDV \leq \rho \\ & D \in \text{EDM}^N \end{array} \quad \begin{array}{ll} \underset{\sqrt[D]{D}}{\text{minimize}} & \| -V(\sqrt[D]{D} - H)V \|^2_F \\ \text{subject to} & \text{rank } VDV \leq \rho \\ & \sqrt[D]{D} \in \sqrt{\text{EDM}^N} \end{array}$$

We apply a convex iteration method for constraining rank. Known heuristics for rank minimization are also explained. We offer new geometrical proof, in §7.1.4.0.1, of a famous discovery by Eckart & Young in 1936 [147]: Euclidean projection on that generally nonconvex subset of the positive semidefinite cone boundary comprising all semidefinite matrices having rank not exceeding a prescribed bound  $\rho$ . We explain how this problem is transformed to a convex optimization for any rank  $\rho$ .

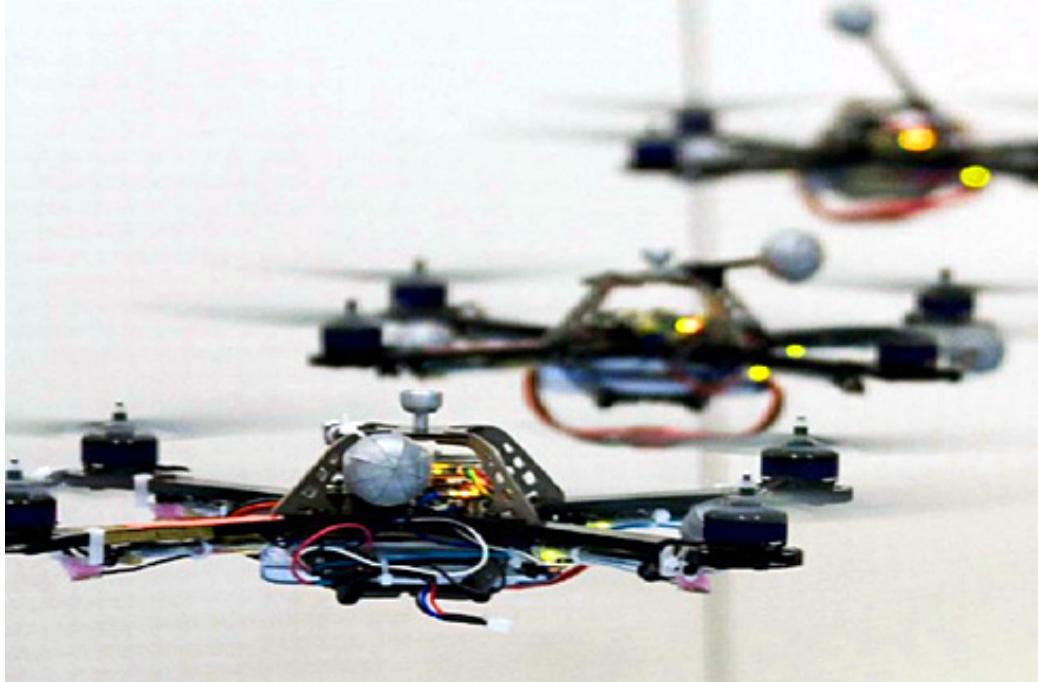


Figure 10: Nanocopter swarm. Robotic vehicles in concert can move larger objects or localize a plume of gas, liquid, or radio waves. [156]

## appendices

We presume a reader already comfortable with elementary vector operations; [14, §3] formally known as *analytic geometry*. [442] Toolboxes are provided, in the form of appendices and code, so as to be more self-contained:

- linear algebra ([appendix A](#) is primarily concerned with proper statements of semidefiniteness for square matrices),
- simple matrices (dyad, doublet, elementary, Householder, Schoenberg, orthogonal, *etcetera*, in [appendix B](#)),
- collection of known analytical solutions to some important optimization problems ([appendix C](#)),
- matrix calculus remains somewhat unsystematized when compared to ordinary calculus ([appendix D](#) concerns matrix-valued functions, matrix differentiation and directional derivatives, Taylor series, and tables of first- and second-order gradients and matrix derivatives),
- an elaborate exposition offering insight into orthogonal and nonorthogonal projection on convex sets (the connection between projection and positive semidefiniteness, for example, or between projection and a linear objective function in [appendix E](#)),
- MATLAB code on [Wikimization](#) [424] to discriminate EDMs, to determine conic independence, to reduce or constrain rank of an optimal solution to a semidefinite program, compressed sensing (compressive sampling) for digital image and audio signal processing, and two distinct methods of reconstructing a map of the United States: one given only distance data, the other given only comparative distance.



Figure 11: Three-dimensional reconstruction of David from distance data.

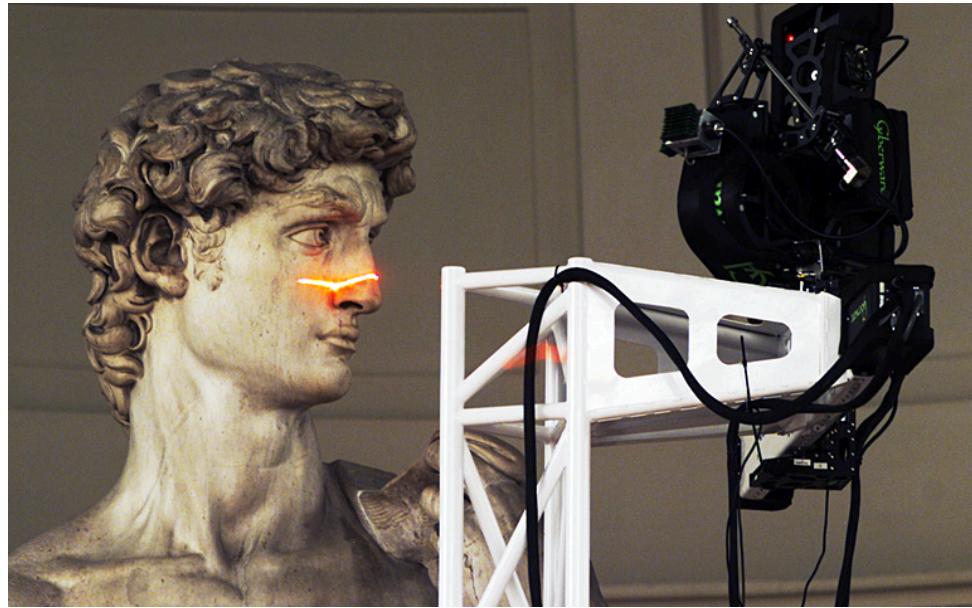


Figure 12: *Digital Michelangelo Project*, Stanford University. Measuring distance to David by laser rangefinder. Spatial resolution is 0.29mm.



# Chapter 2

## Convex geometry

*Convexity has an immensely rich structure and numerous applications. On the other hand, almost every “convex” idea can be explained by a two-dimensional picture.*

— Alexander Barvinok [27, p.vii]

We study convex geometry because it is the easiest of geometries. For that reason, much of a practitioner’s energy is expended seeking invertible transformation of problematic sets to convex ones.

As convex geometry and linear algebra are inextricably bonded by linear inequality (*asymmetry*), we provide much background material on linear algebra (especially in the appendices) although a reader is assumed comfortable with [368] [370] [228] or any other intermediate-level text. The essential references to convex analysis are [225] [343]. The reader is referred to [366] [27] [434] [43] [65] [340] [399] for a comprehensive treatment of convexity. There is relatively less published pertaining to convex matrix-valued functions. [242] [229, §6.6] [329]

### 2.1 Convex set

A set  $\mathcal{C}$  is convex iff for all  $Y, Z \in \mathcal{C}$  and  $0 \leq \mu \leq 1$

$$\mu Y + (1 - \mu)Z \in \mathcal{C} \quad (1)$$

Under that defining condition on  $\mu$ , the linear sum in (1) is called a *convex combination* of  $Y$  and  $Z$ . If  $Y$  and  $Z$  are points in real finite-dimensional Euclidean *vector space* [254] [442]  $\mathbb{R}^n$  or  $\mathbb{R}^{m \times n}$  (matrices), then (1) represents the closed line segment joining them. Line segments are thereby convex sets;  $\mathcal{C}$  is convex iff the line segment connecting any two points in  $\mathcal{C}$  is itself in  $\mathcal{C}$ . Apparent from this definition: a convex set is a connected set. [289, §3.4, §3.5] [43, p.2] A convex set can, but does not necessarily, contain the *origin*  $\mathbf{0}$ .

An *ellipsoid* centered at  $x = a$  (Figure 15 p.36), given matrix  $C \in \mathbb{R}^{m \times n}$

$$\{x \in \mathbb{R}^n \mid \|C(x - a)\|^2 = (x - a)^T C^T C (x - a) \leq 1\} \quad (2)$$

is a good icon for a convex set.<sup>2.1</sup>

---

<sup>2.1</sup>Ellipsoid semiaxes are eigenvectors of  $C^T C$  whose lengths are inverse square root eigenvalues. This particular definition is slablike (Figure 13) in  $\mathbb{R}^n$  when  $C$  has nontrivial nullspace.

### 2.1.1 subspace

A nonempty subset  $\mathcal{R}$  of real Euclidean vector space  $\mathbb{R}^n$  is called a *subspace* (§2.5) if every vector<sup>2.2</sup> of the form  $\alpha x + \beta y$ , for  $\alpha, \beta \in \mathbb{R}$ , is in  $\mathcal{R}$  whenever vectors  $x$  and  $y$  are. [280, §2.3] A subspace is a convex set containing the origin, by definition. [343, p.4] Any subspace is therefore open in the sense that it contains no boundary, but closed in the sense [289, §2]

$$\mathcal{R} + \mathcal{R} = \mathcal{R} \quad (3)$$

It is not difficult to show

$$\mathcal{R} = -\mathcal{R} \quad (4)$$

as is true for any subspace  $\mathcal{R}$ , because  $x \in \mathcal{R} \Leftrightarrow -x \in \mathcal{R}$ . Given any  $x \in \mathcal{R}$

$$\mathcal{R} = x + \mathcal{R} \quad (5)$$

Intersection of an arbitrary collection of subspaces remains a subspace. Any subspace, not constituting the entire *ambient vector space*  $\mathbb{R}^n$ , is a *proper subspace*; e.g.,<sup>2.3</sup> any line (of infinite extent) through the origin in two-dimensional Euclidean space  $\mathbb{R}^2$ . Subspace  $\{\mathbf{0}\}$ , comprising only the origin, is proper though *trivial*. The vector space  $\mathbb{R}^n$  is itself a conventional subspace, inclusively, [254, §2.1] although not proper.

### 2.1.2 linear independence

Arbitrary given vectors in Euclidean space  $\{\Gamma_i \in \mathbb{R}^n, i=1 \dots N\}$  are *linearly independent* (l.i.) if and only if, for all  $\zeta \in \mathbb{R}^N$  ( $\zeta_i \in \mathbb{R}$ )

$$\Gamma_1 \zeta_1 + \cdots + \Gamma_{N-1} \zeta_{N-1} - \Gamma_N \zeta_N = \mathbf{0} \quad (6)$$

has only the *trivial solution*  $\zeta = \mathbf{0}$ ; in other words, iff no vector from the given set can be expressed as a linear combination of those remaining.

Geometrically, two nontrivial vector subspaces are linearly independent iff they intersect only at the origin.

#### 2.1.2.1 preservation of linear independence

(confer §2.4.2.4, §2.10.1) Linear transformation preserves linear dependence. [254, p.86] Conversely, linear independence can be preserved under linear transformation. Given  $Y = [y_1 \ y_2 \ \cdots \ y_N] \in \mathbb{R}^{N \times N}$ , consider the mapping

$$T(\Gamma) : \mathbb{R}^{n \times N} \rightarrow \mathbb{R}^{n \times N} \triangleq \Gamma Y \quad (7)$$

whose domain is the set of all matrices  $\Gamma \in \mathbb{R}^{n \times N}$  holding a linearly independent set columnar. Linear independence of  $\{\Gamma y_i \in \mathbb{R}^n, i=1 \dots N\}$  demands, by definition, there exist no nontrivial solution  $\zeta \in \mathbb{R}^N$  to

$$\Gamma y_1 \zeta_1 + \cdots + \Gamma y_{N-1} \zeta_{N-1} - \Gamma y_N \zeta_N = \mathbf{0} \quad (8)$$

By factoring out  $\Gamma$ , we see that triviality is ensured by linear independence of  $\{y_i \in \mathbb{R}^N\}$ .

---

<sup>2.2</sup>A *vector* is assumed, throughout, to be a column vector.

<sup>2.3</sup>We substitute abbreviation *e.g.* in place of the Latin *exempli gratia*; meaning, *for sake of example*.

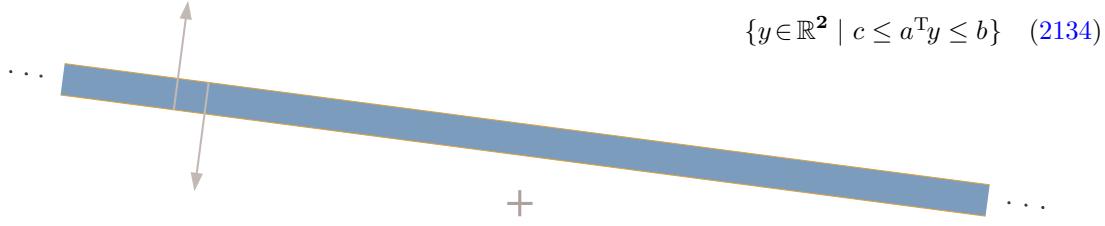


Figure 13: A *slab* is a convex Euclidean body infinite in extent but not affine. Illustrated in  $\mathbb{R}^2$ , it may be constructed by intersecting two opposing halfspaces whose bounding hyperplanes are parallel but not coincident. Because number of halfspaces used in its construction is finite, slab is a *polyhedron* (§2.12). (Cartesian axes + and vector inward-normal, to each halfspace-boundary, are drawn for reference.)

### 2.1.3 Orthant:

name given to a closed convex set that is the higher-dimensional generalization of *quadrant* from the classical Cartesian partition of  $\mathbb{R}^2$ ; a *Cartesian cone*. The most common is the nonnegative orthant  $\mathbb{R}_+^n$  or  $\mathbb{R}_+^{n \times n}$  (analogue to quadrant I) to which membership denotes nonnegative vector- or matrix-entries respectively; *e.g.*,

$$\mathbb{R}_+^n \triangleq \{x \in \mathbb{R}^n \mid x_i \geq 0 \ \forall i\} \quad (9)$$

The nonpositive orthant  $\mathbb{R}_-^n$  or  $\mathbb{R}_-^{n \times n}$  (analogue to quadrant III) denotes negative and 0 entries. Orthant convexity<sup>2.4</sup> is easily verified by definition (1).

### 2.1.4 affine set

A nonempty *affine set* (from the word *affinity*) is any subset of  $\mathbb{R}^n$  that is a translation of some subspace. Any affine set is convex, and open in the sense that it contains no boundary: *e.g.*, empty set  $\emptyset$ , point, line, plane, *hyperplane* (§2.4.2), subspace, *etcetera*. The intersection of an arbitrary collection of affine sets remains affine.

#### 2.1.4.0.1 Definition. Affine subset.

We analogize *affine subset* to subspace,<sup>2.5</sup> defining it to be any nonempty affine set of vectors; an affine subset of  $\mathbb{R}^n$ .  $\triangle$

For some *parallel*<sup>2.6</sup> subspace  $\mathcal{R}$  and any point  $x \in \mathcal{A}$

$$\begin{aligned} \mathcal{A} \text{ is affine} &\Leftrightarrow \mathcal{A} = x + \mathcal{R} \\ &= \{y \mid y - x \in \mathcal{R}\} \end{aligned} \quad (10)$$

*Affine hull* of a set  $\mathcal{C} \subseteq \mathbb{R}^n$  (§2.3.1) is the smallest affine set containing it.

### 2.1.5 dimension

*Dimension* of an arbitrary set  $\mathcal{Z}$  is Euclidean dimension of its affine hull; [434, p.14]

$$\dim \mathcal{Z} \triangleq \dim \text{aff } \mathcal{Z} = \dim \text{aff}(\mathcal{Z} - s), \quad s \in \mathcal{Z} \quad (11)$$

<sup>2.4</sup>All orthants are selfdual simplicial cones. (§2.13.6.1, §2.12.3.1.1)

<sup>2.5</sup>The popular term *affine subspace* is an oxymoron.

<sup>2.6</sup>Two affine sets are *parallel* when one is a translation of the other. [343, p.4]

the same as dimension of the subspace parallel to that affine set  $\text{aff } \mathcal{Z}$  when nonempty. Hence dimension (of a set) is synonymous with *affine dimension*. [225, A.2.1]

### 2.1.6 empty set *versus* empty interior

*Emptiness*  $\emptyset$  of a set is handled differently than *interior* in the classical literature. It is common for a nonempty convex set to have empty interior; *e.g.*, paper in the real world:

- An ordinary flat sheet of paper is a nonempty convex set having empty interior in  $\mathbb{R}^3$  but nonempty interior relative to its affine hull.

#### 2.1.6.1 relative interior

Although it is always possible to pass to a smaller ambient Euclidean space where a nonempty set acquires an interior [27, §II.2.3], we prefer the qualifier *relative* which is the conventional fix to this ambiguous terminology.<sup>2.7</sup> So we distinguish *interior* from *relative interior* throughout: [366] [434] [399]

- Classical interior  $\text{intr } \mathcal{C}$  is defined as a union of points:  $x$  is an interior point of  $\mathcal{C} \subseteq \mathbb{R}^n$  if there exists an open ball of dimension  $n$  and nonzero radius centered at  $x$  that is contained in  $\mathcal{C}$ .
- Relative interior  $\text{rel intr } \mathcal{C}$  of a convex set  $\mathcal{C} \subseteq \mathbb{R}^n$  is interior relative to its affine hull.<sup>2.8</sup>

Thus defined, it is common (though confusing) for  $\text{intr } \mathcal{C}$  the interior of  $\mathcal{C}$  to be empty while its relative interior is not: this happens whenever dimension of its affine hull is less than dimension of the ambient space ( $\dim \text{aff } \mathcal{C} < n$ ; *e.g.*, were  $\mathcal{C}$  paper) or in the exception when  $\mathcal{C}$  is a single point; [289, §2.2.1]

$$\text{rel intr}\{x\} \triangleq \text{aff}\{x\} = \{x\}, \quad \text{intr}\{x\} = \emptyset, \quad x \in \mathbb{R}^n \quad (12)$$

In any case, *closure* of the relative interior of a convex set  $\mathcal{C}$  always yields closure of the set itself;

$$\overline{\text{rel intr } \mathcal{C}} = \overline{\mathcal{C}} \quad (13)$$

Closure is invariant to translation. If  $\mathcal{C}$  is convex then  $\text{rel intr } \mathcal{C}$  and  $\overline{\mathcal{C}}$  are convex. [225, p.24] If  $\mathcal{C}$  has nonempty interior, then

$$\text{rel intr } \mathcal{C} = \text{intr } \mathcal{C} \quad (14)$$

Given the intersection of convex set  $\mathcal{C}$  with affine set  $\mathcal{A}$

$$\text{rel intr}(\mathcal{C} \cap \mathcal{A}) = \text{rel intr}(\mathcal{C}) \cap \mathcal{A} \iff \text{rel intr}(\mathcal{C}) \cap \mathcal{A} \neq \emptyset \quad (15)$$

Because an affine set  $\mathcal{A}$  is open

$$\text{rel intr } \mathcal{A} = \mathcal{A} \quad (16)$$

---

<sup>2.7</sup>Superfluous mingling of terms as in *relatively nonempty set* would be an unfortunate consequence. From the opposite perspective, some authors use the term *full* or *full-dimensional* to describe a set having nonempty interior.

<sup>2.8</sup>Likewise for *relative boundary* (§2.1.7.2), although *relative closure* is superfluous. [225, §A.2.1]

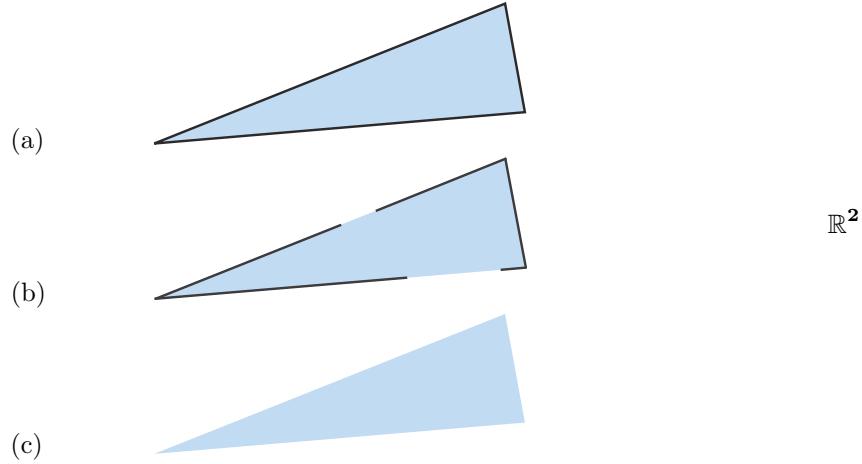


Figure 14: (a) Closed convex set. (b) Neither open, closed, or convex. Yet PSD cone can remain convex in absence of certain boundary components (§2.9.2.9.3). Nonnegative orthant with origin excluded (§2.6) and positive orthant with origin adjoined [343, p.49] are convex. (c) Open convex set.

### 2.1.7 classical boundary

(confer §2.1.7.2) *Boundary* of a set  $\mathcal{C}$  is the closure of  $\mathcal{C}$  less its interior;

$$\partial\mathcal{C} = \overline{\mathcal{C}} \setminus \text{intr } \mathcal{C} \quad (17)$$

[58, §1.1] which follows from the fact

$$\overline{\text{intr } \mathcal{C}} = \overline{\mathcal{C}} \Leftrightarrow \partial \text{intr } \mathcal{C} = \partial \mathcal{C} \quad (18)$$

and presumption of nonempty interior.<sup>2.9</sup> Implications are:

- $\text{intr } \mathcal{C} = \overline{\mathcal{C}} \setminus \partial \mathcal{C}$
- a bounded open set has *boundary* defined but not contained in the set
- interior of an open set is equivalent to the set itself;

from which an open set is defined: [289, p.109]

$$\mathcal{C} \text{ is open} \Leftrightarrow \text{intr } \mathcal{C} = \mathcal{C} \quad (19)$$

$$\mathcal{C} \text{ is closed} \Leftrightarrow \overline{\text{intr } \mathcal{C}} = \mathcal{C} \quad (20)$$

The set illustrated in Figure 14b is not open because it is not equivalent to its interior, for example, it is not closed because it does not contain its boundary, and it is not convex because it does not contain all convex combinations of its boundary points.

<sup>2.9</sup>Otherwise, for  $x \in \mathbb{R}^n$  as in (12), [289, §2.1-§2.3]

$$\overline{\text{intr}\{x\}} = \overline{\emptyset} = \emptyset$$

the empty set is both open and closed.

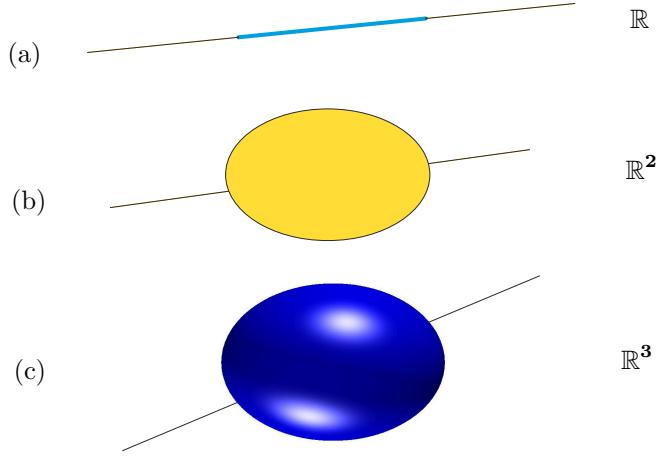


Figure 15: (a) Ellipsoid in  $\mathbb{R}$  is a line segment whose boundary comprises two points. Intersection of line with ellipsoid in  $\mathbb{R}$ , (b) in  $\mathbb{R}^2$ , (c) in  $\mathbb{R}^3$ . Each ellipsoid illustrated has entire boundary constituted by zero-dimensional faces; in fact, by *vertices* (§2.6.1.0.1). Intersection of line with boundary is a point at entry to interior. These same facts hold in higher dimension.

### 2.1.7.1 Line intersection with boundary

A line can intersect the boundary of a convex set in any dimension at a point demarcating the line's entry to the set interior. On one side of that entry-point along the line is the exterior of the set, on the other side is the set interior. In other words,

- starting from any point of a convex set, a move toward the interior is an immediate entry into the interior. [27, §II.2]

When a line intersects the interior of a convex body in any dimension, the boundary appears to the line to be as thin as a point. This is intuitively plausible because, for example, a line intersects the boundary of the ellipsoids in Figure 15 at a point in  $\mathbb{R}$ ,  $\mathbb{R}^2$ , and  $\mathbb{R}^3$ . Such thinness is a remarkable fact when pondering visualization of convex polyhedra (§2.12, §5.14.3) in four Euclidean dimensions, for example, having boundaries constructed from other three-dimensional convex polyhedra called *faces*.

We formally define *face* in (§2.6). For now, we observe the boundary of a convex body to be entirely constituted by all its faces of dimension lower than the body itself. Any face of a convex set is convex. For example: The ellipsoids in Figure 15 have boundaries composed only of zero-dimensional faces. The two-dimensional slab in Figure 13 is an unbounded polyhedron having one-dimensional faces making its boundary. The three-dimensional bounded polyhedron in Figure 22 has zero-, one-, and two-dimensional polygonal faces constituting its boundary.

#### 2.1.7.1.1 Example. Intersection of line with boundary in $\mathbb{R}^6$ .

The convex cone of positive semidefinite matrices  $\mathbb{S}_+^3$  (§2.9), in the ambient subspace of symmetric matrices  $\mathbb{S}^3$  (§2.2.0.1), is a six-dimensional Euclidean body in *isometrically isomorphic*  $\mathbb{R}^6$  (§2.2.1). Boundary of the positive semidefinite cone, in this dimension, comprises faces having only the dimensions 0, 1, and 3; *id est*,  $\{\rho(\rho+1)/2, \rho=0,1,2\}$ .

Unique minimum-distance projection  $PX$  (§E.9) of any point  $X \in \mathbb{S}^3$  on that cone  $\mathbb{S}_+^3$  is known in closed form (§7.1.2). Given, for example,  $\lambda \in \text{intr } \mathbb{R}_+^3$  and *diagonalization* (§A.5.1) of exterior point

$$X = Q\Lambda Q^T \in \mathbb{S}^3, \quad \Lambda \triangleq \begin{bmatrix} \lambda_1 & & \mathbf{0} \\ & \lambda_2 & \\ \mathbf{0} & & -\lambda_3 \end{bmatrix} \quad (21)$$

where  $Q \in \mathbb{R}^{3 \times 3}$  is an *orthogonal matrix*, then the projection on  $\mathbb{S}_+^3$  in  $\mathbb{R}^6$  is

$$PX = Q \begin{bmatrix} \lambda_1 & & \mathbf{0} \\ & \lambda_2 & \\ \mathbf{0} & & 0 \end{bmatrix} Q^T \in \mathbb{S}_+^3 \quad (22)$$

This positive semidefinite matrix  $PX$  nearest  $X$  thus has rank 2, found by discarding all negative eigenvalues in  $\Lambda$ . The line connecting these two points is  $\{X + (PX - X)t \mid t \in \mathbb{R}\}$  where  $t=0 \Leftrightarrow X$  and  $t=1 \Leftrightarrow PX$ . Because this line intersects the boundary of the *positive semidefinite cone*  $\mathbb{S}_+^3$  at point  $PX$  and passes through its interior (by assumption), then the matrix corresponding to an infinitesimally positive perturbation of  $t$  there should reside interior to the cone (rank 3). Indeed, for  $\varepsilon$  an arbitrarily small positive constant,

$$X + (PX - X)t|_{t=1+\varepsilon} = Q(\Lambda + (P\Lambda - \Lambda)(1+\varepsilon))Q^T = Q \begin{bmatrix} \lambda_1 & & \mathbf{0} \\ & \lambda_2 & \\ \mathbf{0} & & \varepsilon\lambda_3 \end{bmatrix} Q^T \in \text{intr } \mathbb{S}_+^3 \quad (23)$$

□

#### 2.1.7.1.2 Example. Tangential line intersection with boundary.

A higher-dimensional boundary  $\partial C$  of a convex Euclidean body  $C$  is simply a dimensionally larger set through which a line can pass when it does not intersect the body's interior. Still, for example, a line existing in five or more dimensions may pass *tangentially* (intersecting no point interior to  $C$  [387, §15.3]) through a single point relatively interior to a three-dimensional face on  $\partial C$ . Let's understand why by inductive reasoning.

Figure 16a shows a vertical line-segment whose boundary comprises its two endpoints. For a line to pass through the boundary tangentially (intersecting no point relatively interior to the line-segment), it must exist in an ambient space of at least two dimensions. Otherwise, the line is confined to the same one-dimensional space as the line-segment and must pass along the segment to reach the end points.

Figure 16b illustrates a two-dimensional ellipsoid whose boundary is constituted entirely by zero-dimensional faces. Again, a line must exist in at least two dimensions to tangentially pass through any single arbitrarily chosen point on the boundary (without intersecting the ellipsoid interior).

Now let's move to an ambient space of three dimensions. Figure 16c shows a polygon rotated into three dimensions. For a line to pass through its zero-dimensional boundary (one of its *vertices*) tangentially, it must exist in at least the two dimensions of the polygon. But for a line to pass tangentially through a single arbitrarily chosen point in the relative interior of a one-dimensional face on the boundary as illustrated, it must exist in at least three dimensions.

Figure 16d illustrates a solid circular cone (drawn truncated) whose one-dimensional faces are halflines emanating from its pointed end (*vertex*). This cone's boundary is constituted solely by those one-dimensional halflines. A line may pass through the boundary tangentially, striking only one arbitrarily chosen point relatively interior to a one-dimensional face, if it exists in at least the three-dimensional ambient space of the cone.

From these few examples, we deduce a general rule (without proof):

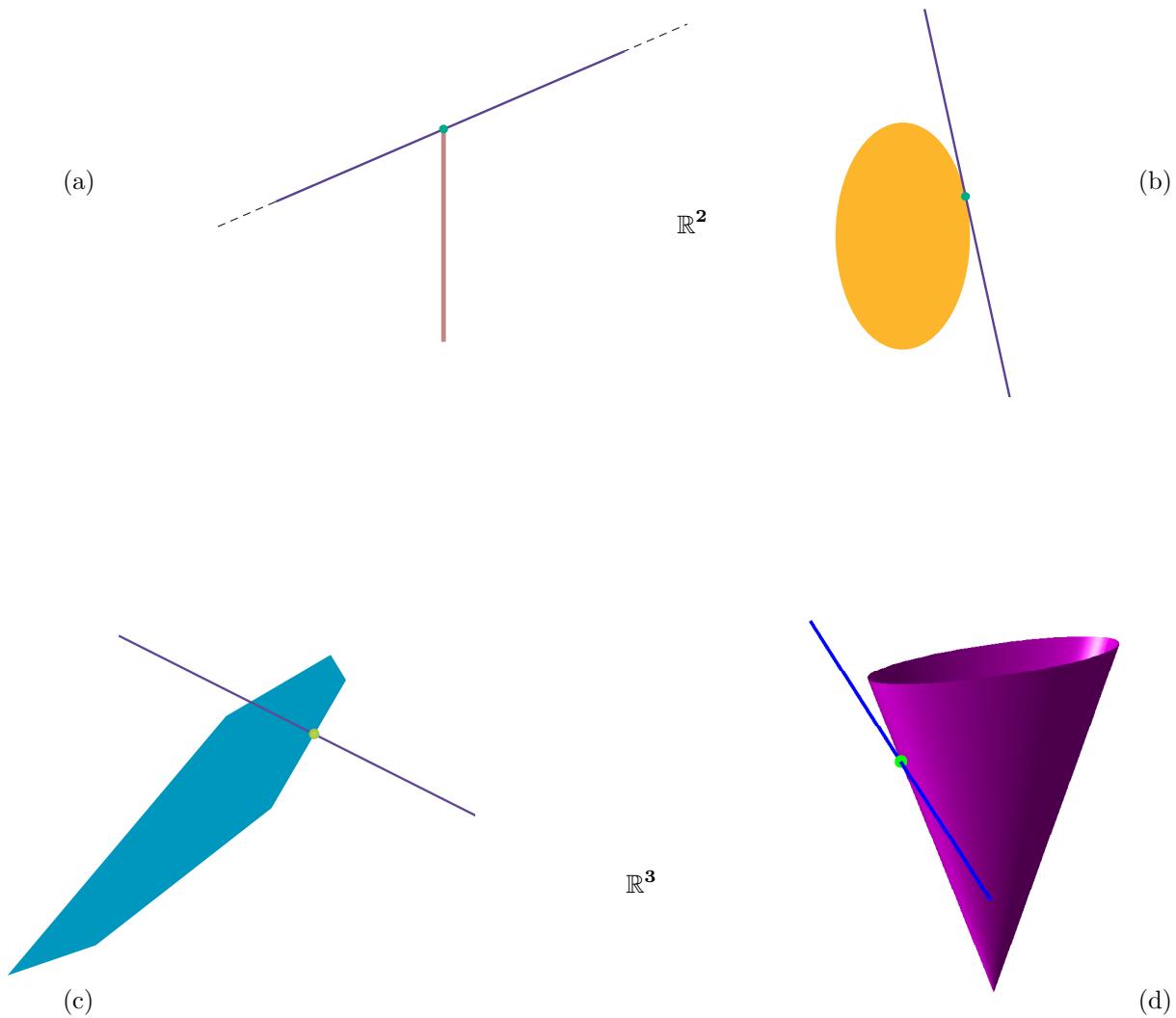


Figure 16: Line tangential: (a) (b) to relative interior of a zero-dimensional face in  $\mathbb{R}^2$ , (c) (d) to relative interior of a one-dimensional face in  $\mathbb{R}^3$ .

- A line may pass tangentially through a single arbitrarily chosen point relatively interior to a  $k$ -dimensional face on the boundary of a convex Euclidean body if the line exists in dimension at least equal to  $k+2$ .

Now the interesting part, with regard to Figure 22 showing a bounded polyhedron in  $\mathbb{R}^3$ ; call it  $\mathcal{P}$ : A line existing in at least four dimensions is required in order to pass tangentially (without hitting  $\text{intr } \mathcal{P}$ ) through a single arbitrary point in the relative interior of any two-dimensional polygonal face on the boundary of polyhedron  $\mathcal{P}$ . Now imagine that polyhedron  $\mathcal{P}$  is itself a three-dimensional face of some other polyhedron in  $\mathbb{R}^4$ . To pass a line tangentially through polyhedron  $\mathcal{P}$  itself, striking only one point from its relative interior  $\text{rel intr } \mathcal{P}$  as claimed, requires a line existing in at least five dimensions.<sup>2.10</sup>

It is not too difficult to deduce:

- A line may pass through a single arbitrarily chosen point interior to a  $k$ -dimensional convex Euclidean body (hitting no other interior point) if that line exists in dimension at least equal to  $k+1$ .

In layman's terms, this means: a being capable of navigating four spatial dimensions (one Euclidean dimension beyond our physical reality) could see inside three-dimensional objects.  $\square$

### 2.1.7.2 Relative boundary

The classical definition of *boundary* of a set  $\mathcal{C}$  presumes nonempty interior:

$$\partial \mathcal{C} = \bar{\mathcal{C}} \setminus \text{intr } \mathcal{C} \quad (17)$$

More suitable to study of convex sets is the *relative boundary*; defined [225, §A.2.1.2]

$$\text{rel } \partial \mathcal{C} \triangleq \bar{\mathcal{C}} \setminus \text{rel intr } \mathcal{C} \quad (24)$$

boundary relative to affine hull of  $\mathcal{C}$ .

In the exception when  $\mathcal{C}$  is a single point  $\{x\}$ , (12)

$$\text{rel } \partial \{x\} = \overline{\{x\}} \setminus \{x\} = \emptyset, \quad x \in \mathbb{R}^n \quad (25)$$

A bounded convex polyhedron (§2.3.2, §2.12.0.0.1) in subspace  $\mathbb{R}$ , for example, has boundary constructed from two points, in  $\mathbb{R}^2$  from at least three line segments, in  $\mathbb{R}^3$  from convex polygons, while a convex *polychoron* (a bounded polyhedron in  $\mathbb{R}^4$  [436]) has boundary constructed from three-dimensional convex polyhedra. A halfspace is partially bounded by a hyperplane; its interior therefore excludes that hyperplane. An affine set has no relative boundary.

## 2.1.8 intersection, sum, difference, product

### 2.1.8.0.1 Theorem. Intersection.

[343, §2, thm.6.5]

Intersection of an arbitrary collection of convex sets  $\{\mathcal{C}_i\}$  is convex. For a finite collection of  $N$  sets, a necessarily nonempty intersection of relative interior  $\bigcap_{i=1}^N \text{rel intr } \mathcal{C}_i = \text{rel intr } \bigcap_{i=1}^N \mathcal{C}_i$  equals relative interior of intersection. And for a possibly infinite collection,  $\bigcap \bar{\mathcal{C}}_i = \overline{\bigcap \mathcal{C}_i}$ .  $\diamond$

In converse this theorem is implicitly false insofar as a convex set can be formed by the intersection of sets that are not. Unions of convex sets are generally not convex. [225, p.22]

---

<sup>2.10</sup>This rule can help determine whether there exists unique solution to a convex optimization problem whose *feasible set* is an intersecting line; e.g., the *trilateration* problem (§5.4.2.2.8).

*Vector sum* of two convex sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$  is convex [225, p.24] (a.k.a *Minkowski sum*)

$$\mathcal{C}_1 + \mathcal{C}_2 = \{x + y \mid x \in \mathcal{C}_1, y \in \mathcal{C}_2\} \quad (26)$$

but not necessarily closed unless at least one set is closed and bounded.

By additive inverse, we can similarly define *vector difference* of two convex sets

$$\mathcal{C}_1 - \mathcal{C}_2 = \{x - y \mid x \in \mathcal{C}_1, y \in \mathcal{C}_2\} \quad (27)$$

which is convex. Applying this definition to nonempty convex set  $\mathcal{C}_1$ , its selfdifference  $\mathcal{C}_1 - \mathcal{C}_1$  is generally nonempty, nontrivial, and convex; e.g, for any *convex cone*  $\mathcal{K}$ , (§2.7.2) the set  $\mathcal{K} - \mathcal{K}$  constitutes its affine hull. [343, p.15]

*Cartesian product* of convex sets

$$\mathcal{C}_1 \times \mathcal{C}_2 = \left\{ \begin{bmatrix} x \\ y \end{bmatrix} \mid x \in \mathcal{C}_1, y \in \mathcal{C}_2 \right\} = \begin{bmatrix} \mathcal{C}_1 \\ \mathcal{C}_2 \end{bmatrix} \quad (28)$$

remains convex. The converse also holds; *id est*, a Cartesian product is convex iff each set is. [225, p.23]

Convex results are also obtained for scaling  $\kappa \mathcal{C}$  of a convex set  $\mathcal{C}$ , rotation/reflection  $Q\mathcal{C}$ , or translation  $\mathcal{C} + \alpha$ ; each similarly defined.

Given any operator  $T$  and convex set  $\mathcal{C}$ , we are prone to write  $T(\mathcal{C})$  meaning

$$T(\mathcal{C}) \triangleq \{T(x) \mid x \in \mathcal{C}\} \quad (29)$$

Given linear operator  $T$ , it therefore follows from (26),

$$\begin{aligned} T(\mathcal{C}_1 + \mathcal{C}_2) &= \{T(x + y) \mid x \in \mathcal{C}_1, y \in \mathcal{C}_2\} \\ &= \{T(x) + T(y) \mid x \in \mathcal{C}_1, y \in \mathcal{C}_2\} \\ &= T(\mathcal{C}_1) + T(\mathcal{C}_2) \end{aligned} \quad (30)$$

### 2.1.9 inverse image

While *epigraph* (§3.5) of a convex function must be convex, it generally holds that inverse image (Figure 17) of a convex function is not. The most prominent examples to the contrary are affine functions (§3.4):

#### 2.1.9.0.1 Theorem. Inverse image.

[343, §3]

Let  $f$  be a mapping from  $\mathbb{R}^{p \times k}$  to  $\mathbb{R}^{m \times n}$ .

- The image of a convex set  $\mathcal{C}$  under any affine function

$$f(\mathcal{C}) = \{f(X) \mid X \in \mathcal{C}\} \subseteq \mathbb{R}^{m \times n} \quad (31)$$

is convex.

- Inverse image of a convex set  $\mathcal{F}$ ,

$$f^{-1}(\mathcal{F}) = \{X \mid f(X) \in \mathcal{F}\} \subseteq \mathbb{R}^{p \times k} \quad (32)$$

a single- or many-valued mapping, under any affine function  $f$  is convex.  $\diamond$

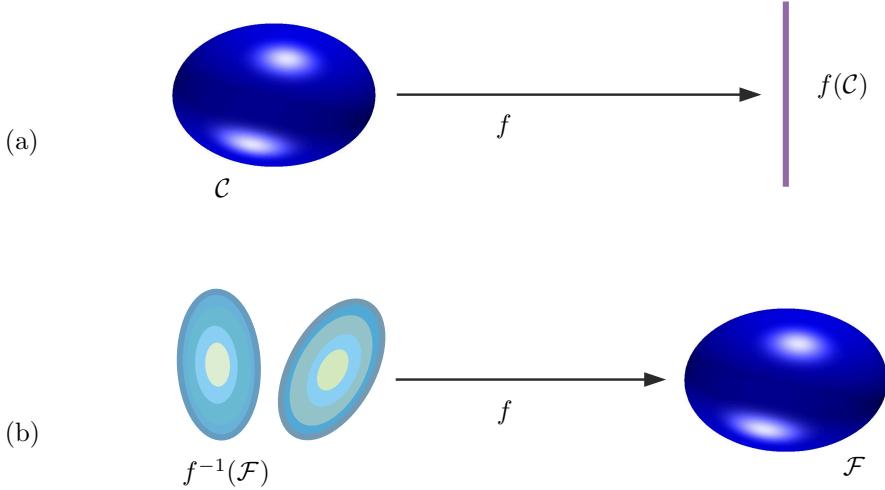


Figure 17: (a) Image of convex set in domain of any convex function  $f$  is convex, but there is no converse. (b) Inverse image under convex function  $f$ .

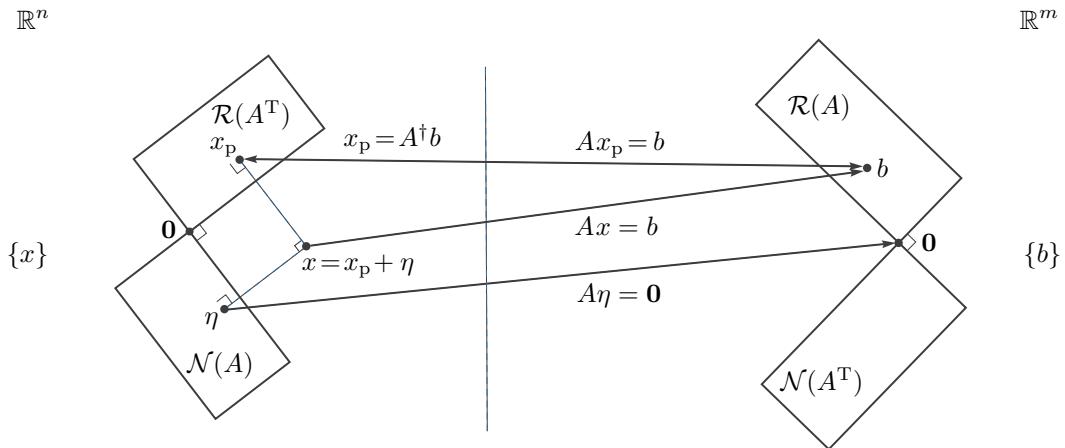


Figure 18: (confer Figure 187) Action of linear map represented by  $A \in \mathbb{R}^{m \times n}$ : [368, p.140] Component of vector  $x$  in nullspace  $\mathcal{N}(A)$  maps to origin while component in rowspace  $\mathcal{R}(A^T)$  maps to range  $\mathcal{R}(A)$ . For any  $A \in \mathbb{R}^{m \times n}$ ,  $A^\dagger Ax = x_p$  and  $AA^\dagger Ax = b$  (§E) and inverse image of  $b \in \mathcal{R}(A)$  is a nonempty affine set:  $x_p + \mathcal{N}(A)$ .

In particular, any affine transformation of an affine set remains affine. [343, p.8] Inverse of any affine transformation, whose image is nonempty and affine, is affine. [343, p.7] Ellipsoids are invariant to any [*sic*] affine transformation.

Although not precluded, this *inverse image theorem* does not require a uniquely invertible mapping  $f$ . Figure 18, for example, mechanizes inverse image under a general linear map. Example 2.9.1.0.2 and §3.5.1 offer further applications.

Each converse of this two-part theorem is generally false; *id est*, given  $f$  affine, a convex image  $f(\mathcal{C})$  does not imply that set  $\mathcal{C}$  is convex, and neither does a convex inverse image  $f^{-1}(\mathcal{F})$  imply set  $\mathcal{F}$  is convex. A counterexample, invalidating a converse, is easy to visualize when the affine function is an orthogonal projector [368] [280]:

**2.1.9.0.2 Corollary.** *Projection on subspace.* [2.11](#) (2118) [343, §3]  
Orthogonal projection of a convex set on a subspace or nonempty affine set is another convex set.  $\diamond$

Again, the converse is false. Shadows, for example, are umbral projections that can be convex when a body providing the shade is not.

## 2.2 Vectorized-matrix inner product

Euclidean space  $\mathbb{R}^n$  comes equipped with a vector inner-product (1061)

$$\langle y, z \rangle \triangleq y^T z = \|y\| \|z\| \cos \psi \quad (33)$$

where  $\psi$  represents angle (in radians) between vectors  $y$  and  $z$ . We prefer those angle brackets to connote a geometric rather than algebraic perspective; *e.g.*, vector  $y$  might represent a hyperplane normal (§2.4.2). Two vectors are *orthogonal (perpendicular)* to one another if and only if their inner product vanishes (iff  $\psi$  is an odd multiple of  $\frac{\pi}{2}$ );

$$y \perp z \Leftrightarrow \langle y, z \rangle = 0 \quad (34)$$

When orthogonal vectors each have unit *norm*, then they are *orthonormal*. A vector inner-product defines Euclidean norm (vector 2-norm, §A.7.1)

$$\|y\|_2 = \|y\| \triangleq \sqrt{y^T y}, \quad \|y\| = 0 \Leftrightarrow y = \mathbf{0} \quad (35)$$

For linear operator  $A$ , its *adjoint*  $A^T$  is a linear operator defined by [254, §3.10]

$$\langle y, A^T z \rangle \triangleq \langle Ay, z \rangle \quad (36)$$

For linear operation on a vector, represented by real matrix  $A$ , the adjoint operator  $A^T$  is its *transposition*. This operator is *selfadjoint* when  $A = A^T$ .

Vector inner-product for matrices is calculated just as it is for vectors; by first transforming a matrix in  $\mathbb{R}^{p \times k}$  to a vector in  $\mathbb{R}^{pk}$  by concatenating its columns in the *natural order*. For lack of a better term, we shall call that linear *bijective* (one-to-one and onto [254, App.A1.2]) transformation *vectorization*. For example, the vectorization of  $Y = [y_1 \ y_2 \ \cdots \ y_k] \in \mathbb{R}^{p \times k}$  [190] [363] is

$$\text{vec } Y \triangleq \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} \in \mathbb{R}^{pk} \quad (37)$$

---

**2.11** For hyperplane representations see §2.4.2. For projection of convex sets on hyperplanes see [434, §6.6]. A nonempty affine set is called an *affine subset* (§2.1.4.0.1). Orthogonal projection of points on affine subsets is reviewed in §E.4.

Then the vectorized-matrix inner-product is trace of matrix inner-product; for  $Z \in \mathbb{R}^{p \times k}$ , [65, §2.6.1] [225, §0.3.1] [444, §8] [406, §2.2]

$$\langle Y, Z \rangle \triangleq \text{tr}(Y^T Z) = \text{vec}(Y)^T \text{vec } Z \quad (38)$$

where (§A.1.1)

$$\text{tr}(Y^T Z) = \text{tr}(Z Y^T) = \text{tr}(Y Z^T) = \text{tr}(Z^T Y) = \mathbf{1}^T (Y \circ Z) \mathbf{1} \quad (39)$$

and where  $\circ$  denotes the *Hadamard product* 2.12 of matrices [181, §1.1.4]. The adjoint  $A^T$  operation on a matrix can therefore be defined in like manner:

$$\langle Y, A^T Z \rangle \triangleq \langle AY, Z \rangle \quad (40)$$

Take any element  $\mathcal{C}_1$  from a matrix-valued set  $\mathcal{C}$  in  $\mathbb{R}^{p \times k}$ , for example, and consider any particular dimensionally compatible real vectors  $v$  and  $w$ . Then vector inner-product of  $\mathcal{C}_1$  with  $vw^T$  is

$$\langle vw^T, \mathcal{C}_1 \rangle = \langle v, \mathcal{C}_1 w \rangle = \langle w^T, v^T \mathcal{C}_1 \rangle = v^T \mathcal{C}_1 w = w^T \mathcal{C}_1^T v = \text{tr}(vw^T \mathcal{C}_1) = \mathbf{1}^T ((vw^T) \circ \mathcal{C}_1) \mathbf{1} \quad (41)$$

Further, linear bijective vectorization is *distributive* with respect to Hadamard product; *id est*,

$$\text{vec}(Y \circ Z) = \text{vec}(Y) \circ \text{vec}(Z) \quad (42)$$

#### 2.2.0.0.1 Example. Application of inverse image theorem.

Suppose set  $\mathcal{C} \subseteq \mathbb{R}^{p \times k}$  were convex. Then for any particular vectors  $v \in \mathbb{R}^p$  and  $w \in \mathbb{R}^k$ , the set of vector inner-products

$$\mathcal{Y} \triangleq v^T \mathcal{C} w = \langle vw^T, \mathcal{C} \rangle \subseteq \mathbb{R} \quad (43)$$

is convex. It is easy to show directly that convex combination of elements from  $\mathcal{Y}$  remains an element of  $\mathcal{Y}$ . 2.13 Instead given convex set  $\mathcal{Y}$ ,  $\mathcal{C}$  must be convex consequent to *inverse image theorem* 2.1.9.0.1.

More generally,  $vw^T$  in (43) may be replaced with any particular matrix  $Z \in \mathbb{R}^{p \times k}$  while convexity of set  $\langle Z, \mathcal{C} \rangle \subseteq \mathbb{R}$  persists. Further, by replacing  $v$  and  $w$  with any particular respective matrices  $U$  and  $W$  of dimension compatible with all elements of convex set  $\mathcal{C}$ , then set  $U^T \mathcal{C} W$  is convex by the *inverse image theorem* because it is a linear mapping of  $\mathcal{C}$ .  $\square$

### 2.2.1 Frobenius'

#### 2.2.1.0.1 Definition. Isomorphic.

An *isomorphism* of a vector space is a transformation equivalent to a linear bijective mapping. Image and inverse image under the transformation operator are then called *isomorphic vector spaces*.  $\triangle$

2.12 Hadamard product is a simple entrywise product of corresponding entries from two matrices of like size; *id est*, not necessarily square. A commutative operation, the Hadamard product can be extracted from within a Kronecker product. [228, p.475]

2.13 To verify that, take any two elements  $\mathcal{C}_1$  and  $\mathcal{C}_2$  from the convex matrix-valued set  $\mathcal{C}$ , and then form the vector inner-products (43) that are two elements of  $\mathcal{Y}$  by definition. Now make a convex combination of those inner products; *videlicet*, for  $0 \leq \mu \leq 1$

$$\mu \langle vw^T, \mathcal{C}_1 \rangle + (1 - \mu) \langle vw^T, \mathcal{C}_2 \rangle = \langle vw^T, \mu \mathcal{C}_1 + (1 - \mu) \mathcal{C}_2 \rangle$$

The two sides are equivalent by linearity of inner product. The right-hand side remains a vector inner-product of  $vw^T$  with an element  $\mu \mathcal{C}_1 + (1 - \mu) \mathcal{C}_2$  from the convex set  $\mathcal{C}$ ; hence, it belongs to  $\mathcal{Y}$ . Since that holds true for any two elements from  $\mathcal{Y}$ , then it must be a convex set.  $\blacklozenge$

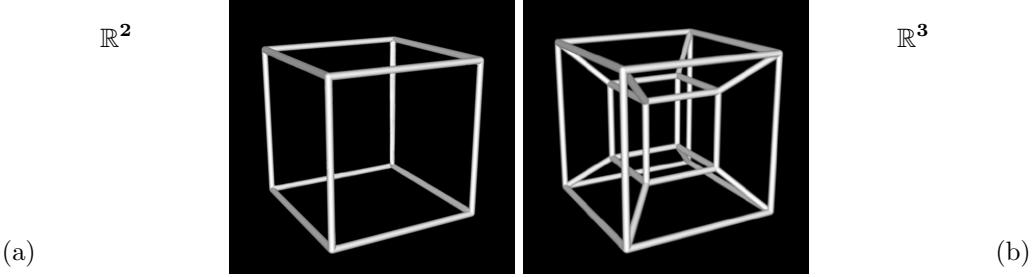


Figure 19: (a) Cube in  $\mathbb{R}^3$  projected on paper-plane  $\mathbb{R}^2$ . Subspace projection operator is not an isomorphism because new adjacencies are introduced. (b) Tesseract is a projection of hypercube in  $\mathbb{R}^4$  on  $\mathbb{R}^3$ .

Isomorphic vector spaces are characterized by preservation of *adjacency*; *id est*, if  $v$  and  $w$  are points connected by a line segment in one vector space, then their images will be connected by a line segment in the other. Two Euclidean bodies may be considered isomorphic if there exists an isomorphism, of their vector spaces, under which the bodies correspond. [408, §I.1] Projection (§E) is not an isomorphism, Figure 19 for example; hence, perfect reconstruction (inverse projection) is generally impossible without additional information.

When  $Z = Y \in \mathbb{R}^{p \times k}$  in (38), *Frobenius' norm* is resultant from vector inner-product; (*confer*(1852))

$$\begin{aligned}\|Y\|_F^2 &= \|\text{vec } Y\|_2^2 = \langle Y, Y \rangle = \text{tr}(Y^T Y) \\ &= \sum_{i,j} Y_{ij}^2 = \sum_i \lambda(Y^T Y)_i = \sum_i \sigma(Y)_i^2\end{aligned}\tag{44}$$

where  $\lambda(Y^T Y)_i$  is the  $i^{\text{th}}$  eigenvalue of  $Y^T Y$ , and  $\sigma(Y)_i$  the  $i^{\text{th}}$  singular value of  $Y$ . Were  $Y$  a *normal matrix* (§A.5.1), then  $\sigma(Y) = |\lambda(Y)|$  [455, §8.1] thus

$$\|Y\|_F^2 = \sum_i \lambda(Y)_i^2 = \|\lambda(Y)\|_2^2 = \langle \lambda(Y), \lambda(Y) \rangle = \langle Y, Y \rangle\tag{45}$$

The converse also holds: [228, §2.5.4]

$$(45) \Rightarrow \text{normal matrix } Y\tag{46}$$

Frobenius' norm is the Euclidean norm of vectorized matrices. Because the metrics are equivalent, for  $X \in \mathbb{R}^{p \times k}$

$$\|\text{vec } X - \text{vec } Y\|_2 = \|X - Y\|_F\tag{47}$$

and because vectorization (37) is a linear bijective map, then vector space  $\mathbb{R}^{p \times k}$  is isometrically isomorphic with vector space  $\mathbb{R}^{pk}$  in the Euclidean sense and  $\text{vec}$  is an *isometric isomorphism* of  $\mathbb{R}^{p \times k}$ . Because of this Euclidean structure, all known results from convex analysis in Euclidean space  $\mathbb{R}^n$  carry over directly to the space of real matrices  $\mathbb{R}^{p \times k}$ ; *e.g.*, norm function convexity (§3.2).

### 2.2.1.1 Injective linear operators

*Injective* mapping (transformation) means one-to-one mapping; synonymous with *uniquely invertible* linear mapping on Euclidean space.

- Linear injective mappings are fully characterized by lack of nontrivial nullspace.

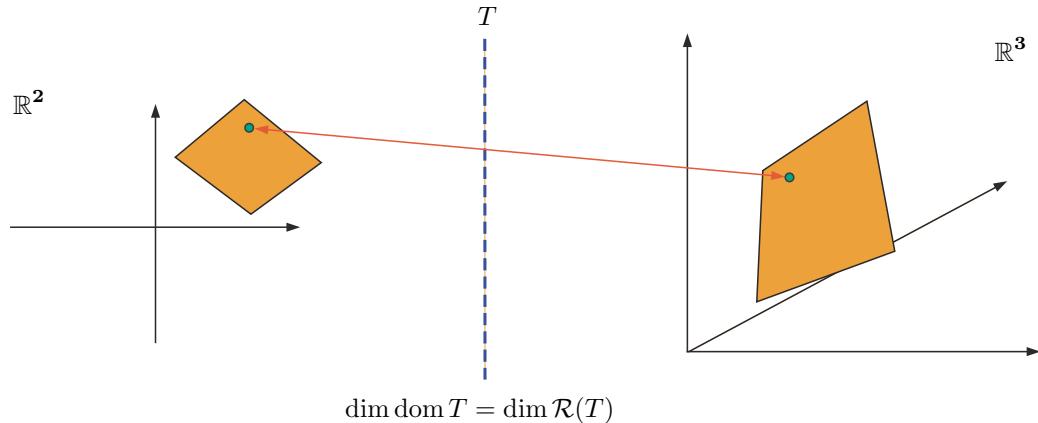


Figure 20: Linear injective mapping  $Tx = Ax : \mathbb{R}^2 \rightarrow \mathbb{R}^3$  of Euclidean body remains two-dimensional under mapping represented by thin full-rank matrix  $A \in \mathbb{R}^{3 \times 2}$ ; two bodies are isomorphic by Definition 2.2.1.0.1.

#### 2.2.1.1.1 Definition. Isometric isomorphism.

An isometric isomorphism of a vector space, having a *metric* defined on it, is a linear bijective mapping  $T$  that preserves distance; *id est*, for all  $x, y \in \text{dom } T$

$$\|Tx - Ty\| = \|x - y\| \quad (48)$$

Then isometric isomorphism  $T$  is called a *bijective isometry*.  $\triangle$

*Unitary linear operator*  $Q : \mathbb{R}^k \rightarrow \mathbb{R}^k$ , represented by orthogonal matrix  $Q \in \mathbb{R}^{k \times k}$  (§B.5.2), is an isometric isomorphism; *e.g.*, discrete Fourier transform via (916). Suppose  $T(X) = UXQ$  is a bijective isometry where  $U$  is a dimensionally compatible *orthonormal matrix*.<sup>2.14</sup> Then we also say Frobenius' norm is *orthogonally invariant*; meaning, for  $X, Y \in \mathbb{R}^{p \times k}$

$$\|U(X - Y)Q\|_F = \|X - Y\|_F \quad (49)$$

Yet isometric operator  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ , represented by  $A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$  on  $\mathbb{R}^2$ , is injective

but not a *surjective* map to  $\mathbb{R}^3$ . [254, §1.6, §2.6] This operator  $T$  can therefore be a bijective isometry only with respect to its range.

Any linear injective transformation on Euclidean space is uniquely invertible on its range. In fact, any linear injective transformation has a range whose dimension equals that of its domain. In other words, for any invertible linear transformation  $T$  [*ibidem*]

$$\dim \text{dom}(T) = \dim \mathcal{R}(T) \quad (50)$$

*e.g.*,  $T$  represented by thin-or-square full-rank matrices. (Figure 20) An important consequence of this fact is:

- Affine dimension, of any  $n$ -dimensional Euclidean body in domain of operator  $T$ , is invariant to linear injective transformation.

---

<sup>2.14</sup> any matrix  $U$  whose columns are orthonormal with respect to each other ( $U^T U = I$ ); these include the orthogonal matrices.

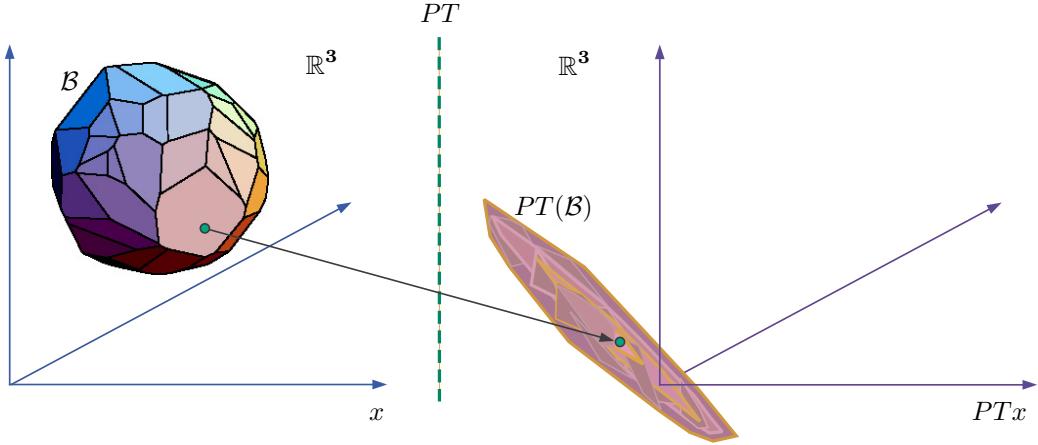


Figure 21: Linear noninjective mapping  $PTx = A^\dagger Ax : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  of three-dimensional Euclidean body  $\mathcal{B}$  has affine dimension 2 under projection on rowspace of wide full-rank matrix  $A \in \mathbb{R}^{2 \times 3}$ . Set of coefficients of orthogonal projection  $T\mathcal{B} = \{Ax \mid x \in \mathcal{B}\}$  is isomorphic with projection  $P(T\mathcal{B})$  [sic].

#### 2.2.1.1.2 Example. Noninjective linear operators.

Mappings in Euclidean space created by noninjective linear operators can be characterized in terms of an orthogonal projector (§E). Consider noninjective linear operator  $Tx = Ax : \mathbb{R}^n \rightarrow \mathbb{R}^m$  represented by wide matrix  $A \in \mathbb{R}^{m \times n}$  ( $m < n$ ). What can be said about the nature of this  $m$ -dimensional mapping?

Concurrently, consider injective linear operator  $Py = A^\dagger y : \mathbb{R}^m \rightarrow \mathbb{R}^n$  where  $\mathcal{R}(A^\dagger) = \mathcal{R}(A^T)$ .  $P(Ax) = PTx$  achieves projection of vector  $x$  on rowspace  $\mathcal{R}(A^T)$ . (§E.3.1) This means vector  $Ax$  can be succinctly interpreted as coefficients of orthogonal projection.

Pseudoinverse matrix  $A^\dagger$  is thin and full-rank, so operator  $Py$  is a linear *bijection* with respect to its range  $\mathcal{R}(A^\dagger)$ . By Definition 2.2.1.0.1, image  $P(T\mathcal{B})$  of projection  $PT(\mathcal{B})$  on  $\mathcal{R}(A^T)$  in  $\mathbb{R}^n$  must therefore be isomorphic with the set of projection coefficients  $T\mathcal{B} = \{Ax \mid x \in \mathcal{B}\}$  in  $\mathbb{R}^m$  and have the same affine dimension by (50). To illustrate, we present a three-dimensional Euclidean body  $\mathcal{B}$  in Figure 21 where any point  $x$  in the nullspace  $\mathcal{N}(A)$  maps to the origin.  $\square$

## 2.2.2 Symmetric matrices

#### 2.2.2.0.1 Definition. Symmetric matrix subspace.

Define a subspace of  $\mathbb{R}^{M \times M}$ : the convex set of all symmetric  $M \times M$  matrices;

$$\mathbb{S}^M \triangleq \left\{ A \in \mathbb{R}^{M \times M} \mid A = A^T \right\} \subseteq \mathbb{R}^{M \times M} \quad (51)$$

This subspace comprising symmetric matrices  $\mathbb{S}^M$  is isomorphic with the vector space  $\mathbb{R}^{M(M+1)/2}$  whose dimension is the number of free variables in a symmetric  $M \times M$  matrix. The *orthogonal complement* [368] [280] of  $\mathbb{S}^M$  is

$$\mathbb{S}^{M \perp} \triangleq \left\{ A \in \mathbb{R}^{M \times M} \mid A = -A^T \right\} \subset \mathbb{R}^{M \times M} \quad (52)$$

the subspace of *antisymmetric* matrices in  $\mathbb{R}^{M \times M}$ ; *id est*,

$$\mathbb{S}^M \oplus \mathbb{S}^{M \perp} = \mathbb{R}^{M \times M} \quad (53)$$

where unique vector sum  $\oplus$  is defined on page 624.  $\triangle$

All antisymmetric matrices have  $\mathbf{0}$  main diagonal by definition. Any square matrix  $A \in \mathbb{R}^{M \times M}$  can be written as a sum of its symmetric and antisymmetric parts: respectively,

$$A = \frac{1}{2}(A + A^T) + \frac{1}{2}(A - A^T) \quad (54)$$

The symmetric part is orthogonal in  $\mathbb{R}^{M^2}$  to the antisymmetric part; *videlicet*,

$$\text{tr}((A + A^T)(A - A^T)) = 0 \quad (55)$$

In the ambient space of real matrices, the antisymmetric matrix subspace can be described

$$\mathbb{S}^{M \perp} = \left\{ \frac{1}{2}(A - A^T) \mid A \in \mathbb{R}^{M \times M} \right\} \subset \mathbb{R}^{M \times M} \quad (56)$$

because any matrix in  $\mathbb{S}^M$  is orthogonal to any matrix in  $\mathbb{S}^{M \perp}$ . Further confined to the ambient subspace of symmetric matrices,  $\mathbb{S}^{M \perp}$  would become trivial ( $\{\mathbf{0}\}$ ).

### 2.2.2.1 Isomorphism of symmetric matrix subspace

When a matrix is symmetric in  $\mathbb{S}^M$ , we may still employ the vectorization transformation (37) to  $\mathbb{R}^{M^2}$ ;  $\text{vec}$ , an isometric isomorphism. We might instead choose to realize in the lower-dimensional subspace  $\mathbb{R}^{M(M+1)/2}$  by ignoring redundant entries (below the main diagonal) during transformation. Such a realization would remain isomorphic but not isometric. Lack of isometry is a spatial distortion due now to disparity in metric between  $\mathbb{R}^{M^2}$  and  $\mathbb{R}^{M(M+1)/2}$ . To realize isometrically in  $\mathbb{R}^{M(M+1)/2}$ , we must make a correction: For  $Y = [Y_{ij}] \in \mathbb{S}^M$  we take symmetric vectorization [242, §2.2.1]

$$\text{svec } Y \triangleq \begin{bmatrix} Y_{11} \\ \sqrt{2}Y_{12} \\ Y_{22} \\ \sqrt{2}Y_{13} \\ \sqrt{2}Y_{23} \\ Y_{33} \\ \vdots \\ Y_{MM} \end{bmatrix} \in \mathbb{R}^{M(M+1)/2} \quad (57)$$

where all entries off the main diagonal have been scaled. Now for  $Z \in \mathbb{S}^M$

$$\langle Y, Z \rangle \triangleq \text{tr}(Y^T Z) = \text{vec}(Y)^T \text{vec } Z = \mathbf{1}^T (Y \circ Z) \mathbf{1} = \text{svec}(Y)^T \text{svec } Z \quad (58)$$

Then because the metrics become equivalent, for  $X \in \mathbb{S}^M$

$$\|\text{svec } X - \text{svec } Y\|_2 = \|X - Y\|_F \quad (59)$$

and because symmetric vectorization (57) is a linear bijective mapping, then  $\text{svec}$  is an isometric isomorphism of the symmetric matrix subspace. In other words,  $\mathbb{S}^M$  is isometrically isomorphic with  $\mathbb{R}^{M(M+1)/2}$  in the Euclidean sense under transformation  $\text{svec}$ .

The set of all symmetric matrices  $\mathbb{S}^M$  forms a proper subspace in  $\mathbb{R}^{M \times M}$ , so for it there exists a standard orthonormal *basis* in isometrically isomorphic  $\mathbb{R}^{M(M+1)/2}$

$$\{E_{ij} \in \mathbb{S}^M\} = \left\{ \begin{array}{ll} e_i e_i^T, & i = j = 1 \dots M \\ \frac{1}{\sqrt{2}}(e_i e_j^T + e_j e_i^T), & 1 \leq i < j \leq M \end{array} \right\} \quad (60)$$

where  $M(M+1)/2$  standard basis matrices  $E_{ij}$  are formed from the standard basis vectors

$$e_i = \left[ \begin{array}{ll} 1, & i = j \\ 0, & i \neq j \end{array}, \quad j = 1 \dots M \right] \in \mathbb{R}^M \quad (61)$$

Thus we have a basic orthogonal expansion for  $Y \in \mathbb{S}^M$

$$Y = \sum_{j=1}^M \sum_{i=1}^j \langle E_{ij}, Y \rangle E_{ij} \quad (62)$$

whose unique coefficients

$$\langle E_{ij}, Y \rangle = \begin{cases} Y_{ii}, & i = 1 \dots M \\ Y_{ij}\sqrt{2}, & 1 \leq i < j \leq M \end{cases} \quad (63)$$

correspond to entries of the symmetric vectorization (57).

### 2.2.3 Symmetric hollow subspace

#### 2.2.3.0.1 Definition. Hollow subspaces. [393]

Define the *real hollow subspace* of  $\mathbb{R}^{M \times M}$  to be the convex set of all symmetric  $M \times M$  matrices having  $\mathbf{0}$  main diagonal;

$$\mathbb{R}_h^{M \times M} \triangleq \{A \in \mathbb{R}^{M \times M} \mid A = A^T, \delta(A) = \mathbf{0}\} \subset \mathbb{R}^{M \times M} \quad (64)$$

where the main diagonal of  $A \in \mathbb{R}^{M \times M}$  is denoted (§A.1)

$$\delta(A) \in \mathbb{R}^M \quad (1569)$$

Operating on a vector, linear operator  $\delta$  naturally returns a diagonal matrix;  $\delta(\delta(A))$  is a diagonal matrix. Operating recursively on a vector  $\Lambda \in \mathbb{R}^N$  or diagonal matrix  $\Lambda \in \mathbb{S}^N$ , operator  $\delta(\delta(\Lambda))$  returns  $\Lambda$  itself;

$$\delta^2(\Lambda) \equiv \delta(\delta(\Lambda)) = \Lambda \quad (1571)$$

The subspace  $\mathbb{R}_h^{M \times M}$  (64) comprising (real) symmetric hollow matrices is isomorphic with subspace  $\mathbb{R}^{M(M-1)/2}$ ; its orthogonal complement is

$$\mathbb{R}_h^{M \times M \perp} \triangleq \{A \in \mathbb{R}^{M \times M} \mid A = -A^T + 2\delta^2(A)\} \subseteq \mathbb{R}^{M \times M} \quad (65)$$

the subspace of *antisymmetric antihollow* matrices in  $\mathbb{R}^{M \times M}$ ; *id est*,

$$\mathbb{R}_h^{M \times M} \oplus \mathbb{R}_h^{M \times M \perp} = \mathbb{R}^{M \times M} \quad (66)$$

Yet defined instead as a proper subspace of ambient  $\mathbb{S}^M$

$$\begin{aligned} \mathbb{S}_h^M &\triangleq \{A \in \mathbb{S}^M \mid \delta(A) = \mathbf{0}\} \subset \mathbb{S}^M \\ &\equiv \mathbb{R}_h^{M \times M} \end{aligned} \quad (67)$$

the orthogonal complement  $\mathbb{S}_h^{M\perp}$  of *symmetric hollow subspace*  $\mathbb{S}_h^M$  (confer (65))

$$\mathbb{S}_h^{M\perp} \triangleq \left\{ A \in \mathbb{S}^M \mid A = \delta^2(A) \right\} \subseteq \mathbb{S}^M \quad (68)$$

(called *symmetric antihollow subspace*) is simply the subspace of diagonal matrices; *id est*,

$$\mathbb{S}_h^M \oplus \mathbb{S}_h^{M\perp} = \mathbb{S}^M \quad (69)$$

having  $\dim \mathbb{S}_h^M = M(M-1)/2$  and  $\dim \mathbb{S}_h^{M\perp} = M$  in isomorphic  $\mathbb{R}^{M(M+1)/2}$ .  $\triangle$

Any matrix  $A \in \mathbb{R}^{M \times M}$  can be written as a sum of its symmetric hollow and antisymmetric antihollow parts: respectively,

$$A = \left( \frac{1}{2}(A + A^T) - \delta^2(A) \right) + \left( \frac{1}{2}(A - A^T) + \delta^2(A) \right) \quad (70)$$

The symmetric hollow part is orthogonal to the antisymmetric antihollow part in  $\mathbb{R}^{M^2}$ ; *videlicet*,

$$\text{tr} \left( \left( \frac{1}{2}(A + A^T) - \delta^2(A) \right) \left( \frac{1}{2}(A - A^T) + \delta^2(A) \right) \right) = 0 \quad (71)$$

because any matrix in subspace  $\mathbb{R}_h^{M \times M}$  is orthogonal to any matrix in the antisymmetric antihollow subspace

$$\mathbb{R}_h^{M \times M \perp} = \left\{ \frac{1}{2}(A - A^T) + \delta^2(A) \mid A \in \mathbb{R}^{M \times M} \right\} \subseteq \mathbb{R}^{M \times M} \quad (72)$$

of the ambient space of real matrices; which reduces to the diagonal matrices in the ambient space of symmetric matrices

$$\mathbb{S}_h^{M\perp} = \left\{ \delta^2(A) \mid A \in \mathbb{S}^M \right\} = \left\{ \delta(u) \mid u \in \mathbb{R}^M \right\} \subseteq \mathbb{S}^M \quad (73)$$

In anticipation of their utility with Euclidean distance matrices (EDMs) in §5, for symmetric hollow matrices we introduce the linear bijective vectorization `dvec` that is the natural analogue to symmetric matrix vectorization `svec` (57): for  $Y = [Y_{ij}] \in \mathbb{S}_h^M$

$$\text{dvec } Y \triangleq \sqrt{2} \begin{bmatrix} Y_{12} \\ Y_{13} \\ Y_{23} \\ Y_{14} \\ Y_{24} \\ Y_{34} \\ \vdots \\ Y_{M-1,M} \end{bmatrix} \in \mathbb{R}^{M(M-1)/2} \quad (74)$$

Like `svec`, `dvec` is an isometric isomorphism on the symmetric hollow subspace. For  $X \in \mathbb{S}_h^M$

$$\| \text{dvec } X - \text{dvec } Y \|_2 = \| X - Y \|_{\text{F}} \quad (75)$$

The set of all symmetric hollow matrices  $\mathbb{S}_h^M$  forms a proper subspace in  $\mathbb{R}^{M \times M}$ , so for it there must be a standard orthonormal basis in isometrically isomorphic  $\mathbb{R}^{M(M-1)/2}$

$$\{E_{ij} \in \mathbb{S}_h^M\} = \left\{ \frac{1}{\sqrt{2}}(e_i e_j^T + e_j e_i^T), \quad 1 \leq i < j \leq M \right\} \quad (76)$$

where  $M(M-1)/2$  standard basis matrices  $E_{ij}$  are formed from the standard basis vectors  $e_i \in \mathbb{R}^M$ .

The *symmetric hollow majorization corollary* A.1.2.0.2 characterizes eigenvalues of symmetric hollow matrices.

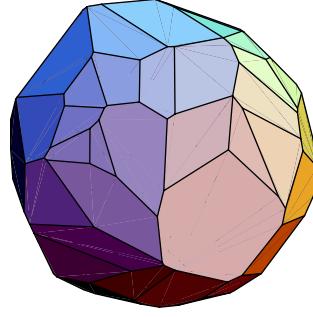


Figure 22: Convex hull of a random list of points in  $\mathbb{R}^3$ . Some points from that generating list reside interior to this *convex polyhedron* (§2.12). [436, *Convex Polyhedron*] (Avis-Fukuda-Mizukoshi)

## 2.3 Hulls

We focus on the affine, convex, and conic hulls: convex sets that may be regarded as kinds of Euclidean container or vessel united with its interior.

### 2.3.1 Affine hull, affine dimension

Affine dimension of any set in  $\mathbb{R}^n$  is the dimension of the smallest affine set (empty set, point, line, plane, hyperplane (§2.4.2), translated subspace,  $\mathbb{R}^n$ ) that contains it. For nonempty sets, affine dimension is the same as dimension of the subspace parallel to that affine set. [343, §1] [225, §A.2.1]

Ascribe the points in a list  $\{x_\ell \in \mathbb{R}^n, \ell=1 \dots N\}$  to the columns of matrix  $X$ :

$$X = [x_1 \cdots x_N] \in \mathbb{R}^{n \times N} \quad (77)$$

In particular, we define *affine dimension*  $r$  of the  $N$ -point list  $X$  as dimension of the smallest affine set in Euclidean space  $\mathbb{R}^n$  that contains  $X$ ;

$$r \triangleq \dim \text{aff } X \quad (78)$$

Affine dimension  $r$  is a lower bound sometimes called *embedding dimension*. [393] [211] That affine set  $\mathcal{A}$  in which those points are embedded is unique and called the *affine hull* [366, §2.1];

$$\begin{aligned} \mathcal{A} &\triangleq \text{aff } \{x_\ell \in \mathbb{R}^n, \ell=1 \dots N\} = \text{aff } X \\ &= x_1 + \mathcal{R}\{x_\ell - x_1, \ell=2 \dots N\} = \{Xa \mid a^T \mathbf{1} = 1\} \subseteq \mathbb{R}^n \end{aligned} \quad (79)$$

for which we call list  $X$  a set of *generators*. Hull  $\mathcal{A}$  is parallel to subspace

$$\mathcal{R}\{x_\ell - x_1, \ell=2 \dots N\} = \mathcal{R}(X - x_1 \mathbf{1}^T) \subseteq \mathbb{R}^n \quad (80)$$

where

$$\mathcal{R}(A) = \{Ax \mid \forall x\} \quad (144)$$

Given some arbitrary set  $\mathcal{C}$  and any  $x \in \mathcal{C}$

$$\text{aff } \mathcal{C} = x + \text{aff}(\mathcal{C} - x) \quad (81)$$

where  $\text{aff}(\mathcal{C} - x)$  is a subspace.

$$\text{aff } \emptyset \triangleq \emptyset \quad (82)$$

The affine hull of a point  $x$  is that point itself;

$$\text{aff}\{x\} = \{x\} \quad (83)$$

Affine hull of two distinct points is the unique line through them. (Figure 23) The affine hull of three noncollinear points in any dimension is that unique plane containing the points, and so on. Affine hull of a convex cone is the same as affine hull of its extreme directions and the origin. The subspace of symmetric matrices  $\mathbb{S}^m$  is the affine hull of the cone of positive semidefinite matrices; (§2.9)

$$\text{aff } \mathbb{S}_+^m = \mathbb{S}^m \quad (84)$$

**2.3.1.0.1 Example.** *Affine hull of rank-1 correlation matrices.* (confer §5.9.1.0.1) [245]  
The set of all  $m \times m$  rank-1 correlation matrices is defined by all binary vectors  $y \in \mathbb{R}^m$

$$\{yy^T \in \mathbb{S}_+^m \mid \delta(yy^T) = \mathbf{1}\} \quad (85)$$

Affine hull of the rank-1 correlation matrices is equal to the set of normalized symmetric matrices; *id est*,

$$\text{aff}\{yy^T \in \mathbb{S}_+^m \mid \delta(yy^T) = \mathbf{1}\} = \{A \in \mathbb{S}^m \mid \delta(A) = \mathbf{1}\} \quad (86)$$

□

**2.3.1.0.2 Exercise.** *Affine hull of correlation matrices.*

Prove (86) via definition of affine hull. Find the convex hull instead. ▼

### 2.3.1.1 Partial order induced by $\mathbb{R}_+^N$ and $\mathbb{S}_+^M$

Notation  $a \succeq 0$  means vector  $a$  belongs to nonnegative orthant  $\mathbb{R}_+^N$  while  $a \succ 0$  means vector  $a$  belongs to the nonnegative orthant's interior  $\text{intr } \mathbb{R}_+^N$ .  $a \succeq b$  denotes comparison of vector  $a$  to vector  $b$  on  $\mathbb{R}^N$  with respect to the nonnegative orthant; *id est*,  $a \succeq b$  means  $a - b$  belongs to the nonnegative orthant but neither  $a$  or  $b$  is necessarily nonnegative. With particular respect to the nonnegative orthant,  $a \succeq b \Leftrightarrow a_i \geq b_i \forall i$  (373).

More generally,  $a \succeq_{\mathcal{K}} b$  denotes comparison with respect to pointed closed convex cone  $\mathcal{K}$ , whereas comparison with respect to the cone's interior is denoted  $a \succ_{\mathcal{K}} b$ . But equivalence with entrywise comparison does not generally hold, and neither  $a$  or  $b$  necessarily belongs to  $\mathcal{K}$ . (§2.7.2.2)

The symbol  $\geq$  is reserved for scalar comparison on the real line  $\mathbb{R}$  with respect to the nonnegative real line  $\mathbb{R}_+$  as in  $a^T y \geq b$ . Comparison of matrices with respect to the positive semidefinite cone  $\mathbb{S}_+^M$ , like  $I \succeq A \succeq 0$  in Example 2.3.2.0.1, is explained in §2.9.0.1.

## 2.3.2 Convex hull

The *convex hull* [225, §A.1.4] [343] of any bounded<sup>2.15</sup> list or set of  $N$  points  $X \in \mathbb{R}^{n \times N}$  forms a unique bounded convex *polyhedron* (confer §2.12.0.0.1) whose vertices constitute

<sup>2.15</sup>An arbitrary set  $\mathcal{C}$  in  $\mathbb{R}^n$  is *bounded* iff it can be contained in a Euclidean ball having finite radius. [126, §2.2] (confer §5.7.3.0.1) The smallest ball containing  $\mathcal{C}$  has radius  $\inf_x \sup_{y \in \mathcal{C}} \|x - y\|$  and center  $x^*$  whose

determination is a convex problem because  $\sup_{y \in \mathcal{C}} \|x - y\|$  is a convex function of  $x$ ; but the supremum may be difficult to ascertain.

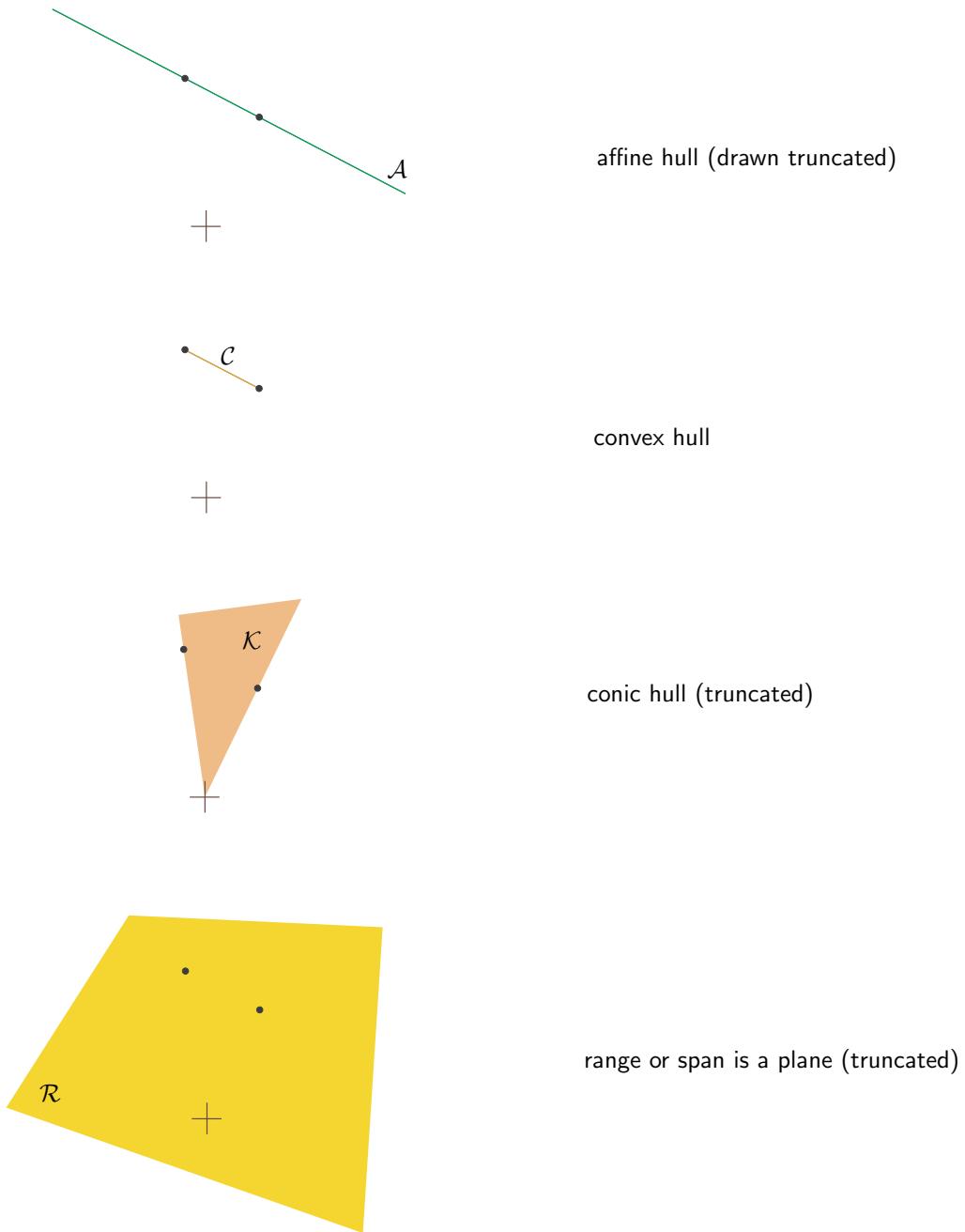


Figure 23: Given two points in Euclidean vector space of any dimension, their various hulls are illustrated. Each hull is a subset of range; generally,  $\mathcal{A}, \mathcal{C}, \mathcal{K} \subseteq \mathcal{R} \ni \mathbf{0}$ . (Cartesian axes drawn for reference.)

some subset of that list;

$$\mathcal{P} \triangleq \text{conv}\{x_\ell, \ell=1\dots N\} = \text{conv } X = \{Xa \mid a^T \mathbf{1} = 1, a \succeq 0\} \subseteq \mathbb{R}^n \quad (87)$$

Union of relative interior and relative boundary (§2.1.7.2) of the polyhedron comprise its convex hull  $\mathcal{P}$ , the smallest closed convex set that contains list  $X$ ; e.g., Figure 22. Given  $\mathcal{P}$ , generating list  $\{x_\ell\}$  is not unique. But because every bounded polyhedron is the convex hull of its vertices, [366, §2.12.2] the vertices of  $\mathcal{P}$  comprise a *minimal set* of generators.

Given some arbitrary set  $\mathcal{C} \subseteq \mathbb{R}^n$ , its convex hull  $\text{conv } \mathcal{C}$  is equivalent to the smallest convex set containing it. (confer §2.4.1.1.1) The convex hull is a subset of the affine hull;

$$\text{conv } \mathcal{C} \subseteq \text{aff } \mathcal{C} = \text{aff } \bar{\mathcal{C}} = \overline{\text{aff } \mathcal{C}} = \text{aff conv } \mathcal{C} \quad (88)$$

Any closed bounded convex set  $\mathcal{C}$  is equal to the convex hull of its boundary;

$$\mathcal{C} = \text{conv } \partial \mathcal{C} \quad (89)$$

$$\text{conv } \emptyset \triangleq \emptyset \quad (90)$$

**2.3.2.0.1 Example.** *Hull of rank- $k$  projection matrices.* [162] [319] [12, §4.1] [326, §3] [268, §2.4] [269] Convex hull of the set comprising outer product of orthonormal matrices has equivalent expression: for  $1 \leq k \leq N$  (§2.9.0.1)

$$\text{conv}\left\{UU^T \mid U \in \mathbb{R}^{N \times k}, U^T U = I\right\} = \left\{A \in \mathbb{S}^N \mid I \succeq A \succeq 0, \langle I, A \rangle = k\right\} \subset \mathbb{S}_+^N \quad (91)$$

This important convex body we call *Fantope* (after mathematician Ky Fan). In case  $k=1$ , there is slight simplification: ((1778), Example 2.9.2.7.1)

$$\text{conv}\left\{UU^T \mid U \in \mathbb{R}^N, U^T U = I\right\} = \left\{A \in \mathbb{S}^N \mid A \succeq 0, \langle I, A \rangle = 1\right\} \quad (92)$$

This particular Fantope is called *spectahedron*. [sic] [170, §5.1] In case  $k=N$ , the Fantope is Identity matrix  $I$ . More generally, the set

$$\left\{UU^T \mid U \in \mathbb{R}^{N \times k}, U^T U = I\right\} \quad (93)$$

comprises the extreme points (§2.6.0.0.1) of its convex hull. By (1617), each and every extreme point  $UU^T$  has only  $k$  nonzero eigenvalues  $\lambda$  and they all equal 1; *id est*,  $\lambda(UU^T)_{1:k} = \lambda(U^T U) = \mathbf{1}$ . So Frobenius' norm of each and every extreme point equals the same constant

$$\|UU^T\|_F^2 = k \quad (94)$$

Each extreme point simultaneously lies on the boundary of the positive semidefinite cone (when  $k < N$ , §2.9) and on the boundary of a *hypersphere* of dimension  $k(N - \frac{k}{2} + \frac{1}{2})$  and radius  $\sqrt{k(1 - \frac{k}{N})}$  centered at  $\frac{k}{N}I$  along the ray (base  $\mathbf{0}$ ) through the Identity matrix  $I$  in isomorphic vector space  $\mathbb{R}^{N(N+1)/2}$  (§2.2.2.1).

Figure 24 illustrates extreme points (93) comprising the boundary of a Fantope: the boundary of a *disc* corresponding to  $k=1$ ,  $N=2$ ; but that circumscription does not hold in higher dimension. (§2.9.2.8)  $\square$

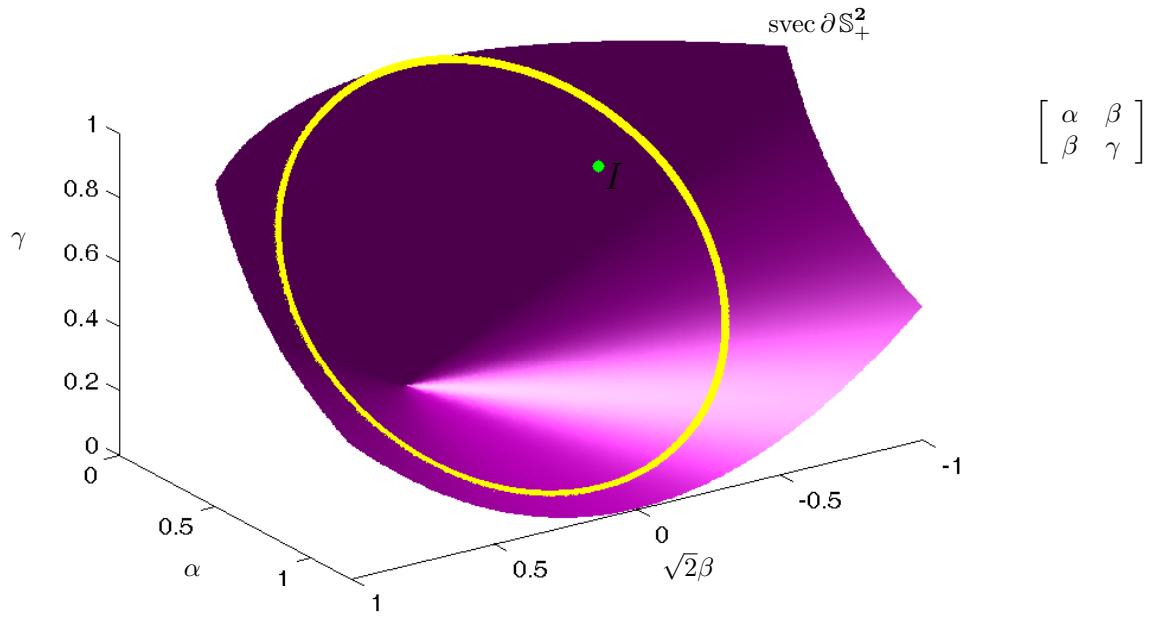


Figure 24: Two Fantopes. Circle (radius  $1/\sqrt{2}$ ), shown here on boundary of positive semidefinite cone  $\mathbb{S}_+^2$  in isometrically isomorphic  $\mathbb{R}^3$  from Figure 46, comprises boundary of a Fantope (91) in this dimension ( $k = 1$ ,  $N = 2$ ). Lone point illustrated is Identity matrix  $I$ , interior to PSD cone, and is that Fantope corresponding to  $k = 2$ ,  $N = 2$ . (View is from inside PSD cone looking toward origin.)

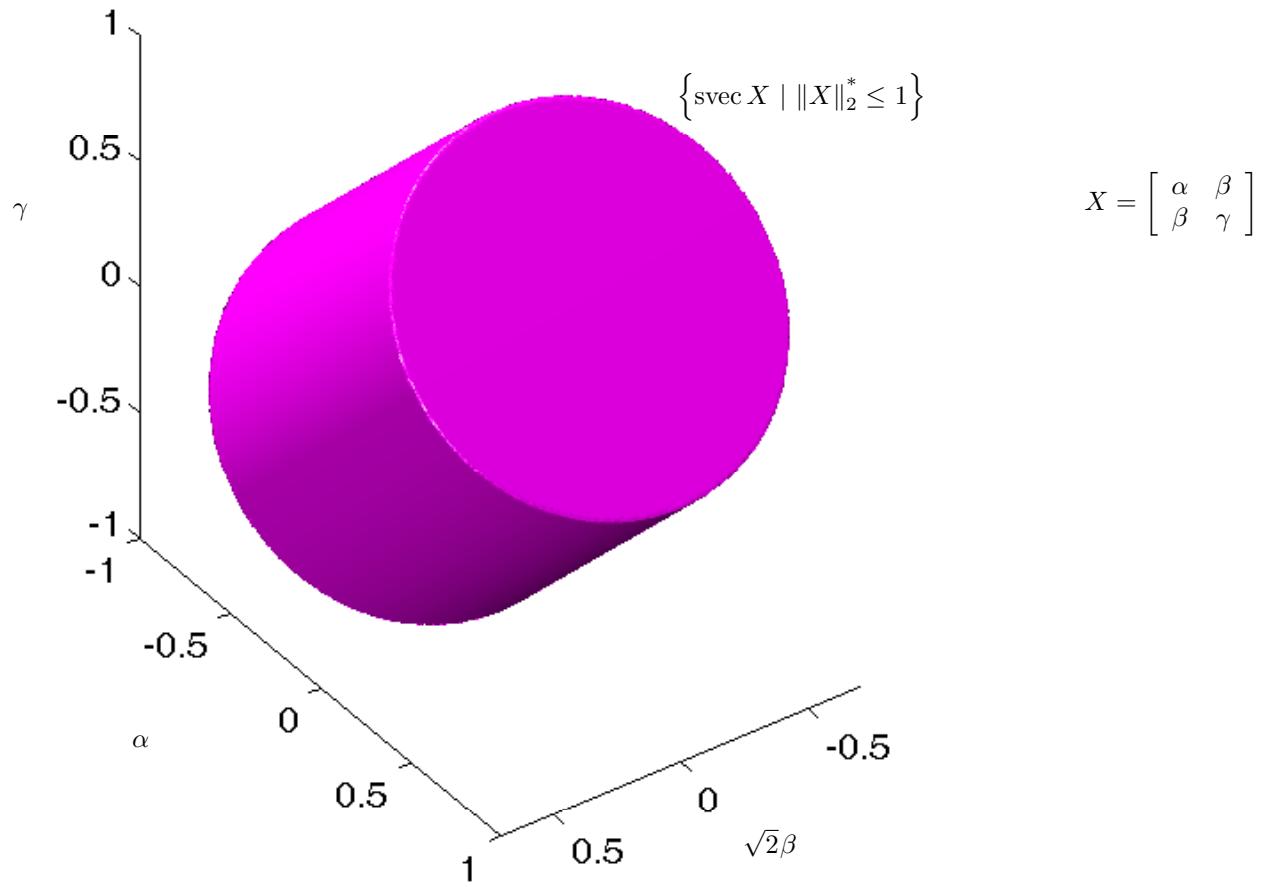


Figure 25: Nuclear norm is a sum of singular values;  $\|X\|_2^* \triangleq \sum_i \sigma(X)_i$ . Nuclear norm ball, in the subspace of  $2 \times 2$  symmetric matrices, is a truncated cylinder in isometrically isomorphic  $\mathbb{R}^3$ .

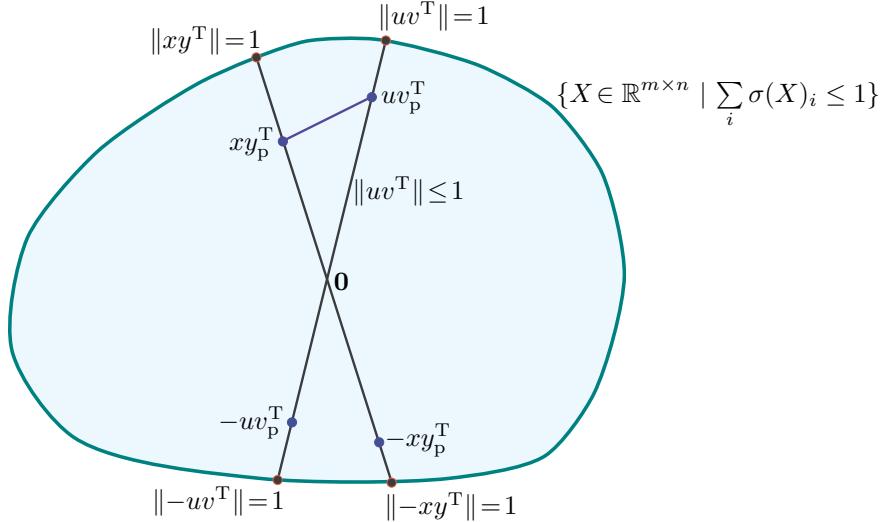


Figure 26:  $uv_p^T$  is a convex combination of normalized dyads  $\|\pm uv^T\|=1$ ; similarly for  $xy_p^T$ . Any point in line segment joining  $xy_p^T$  to  $uv_p^T$  is expressible as a convex combination of two to four points indicated on boundary.

### 2.3.2.0.2 Example. Nuclear norm ball: convex hull of rank-1 matrices.

From (92), in Example 2.3.2.0.1, we learn that the convex hull of normalized symmetric rank-1 matrices is a slice of the positive semidefinite cone. In §2.9.2.7 we find the convex hull of all symmetric rank-1 matrices to be the entire positive semidefinite cone.

In the present example we abandon symmetry; instead posing, what is the convex hull of bounded nonsymmetric rank-1 matrices:

$$\text{conv}\{uv^T \mid \|uv^T\| \leq 1, u \in \mathbb{R}^m, v \in \mathbb{R}^n\} = \{X \in \mathbb{R}^{m \times n} \mid \sum_i \sigma(X)_i \leq 1\} \quad (95)$$

where  $\sigma(X)$  is a vector of singular values. (Since  $\|uv^T\| = \|u\|\|v\|$  (1768), norm of each vector constituting a *dyad*  $uv^T$  (§B.1) in the hull is effectively bounded above by 1.)

**Proof.** ( $\Leftarrow$ ) Suppose  $\sum \sigma(X)_i \leq 1$ . Decompose  $X = U\Sigma V^T$  by SVD (§A.6) where  $U = [u_1 \dots u_{\min\{m,n\}}] \in \mathbb{R}^{m \times \min\{m,n\}}$ ,  $V = [v_1 \dots v_{\min\{m,n\}}] \in \mathbb{R}^{n \times \min\{m,n\}}$ , and whose sum of singular values is  $\sum \sigma(X)_i = \text{tr } \Sigma = \kappa \leq 1$ . Then we may write  $X = \sum \frac{\sigma_i}{\kappa} \sqrt{\kappa} u_i \sqrt{\kappa} v_i^T$  which is a convex combination of dyads each of whose norm does not exceed 1. (Srebro)

( $\Rightarrow$ ) Now suppose we are given a convex combination of dyads  $X = \sum \alpha_i u_i v_i^T$  such that  $\sum \alpha_i = 1$ ,  $\alpha_i \geq 0 \forall i$ , and  $\|u_i v_i^T\| \leq 1 \forall i$ . Then by triangle inequality for singular values [229, cor.3.4.3]  $\sum \sigma(X)_i \leq \sum \sigma(\alpha_i u_i v_i^T) = \sum \alpha_i \|u_i v_i^T\| \leq \sum \alpha_i$ .  $\blacklozenge$

Given any particular dyad  $uv_p^T$  in the convex hull, because its polar  $-uv_p^T$  and every convex combination of the two belong to that hull, then the unique line containing the two points  $\pm uv_p^T$  (their affine combination (79)) must intersect the hull's boundary at the normalized dyads  $\{\pm uv^T \mid \|uv^T\|=1\}$ . Any point formed by convex combination of dyads in the hull must therefore be expressible as a convex combination of dyads on the boundary: Figure 26,

$$\text{conv}\{uv^T \mid \|uv^T\| \leq 1, u \in \mathbb{R}^m, v \in \mathbb{R}^n\} \equiv \text{conv}\{uv^T \mid \|uv^T\| = 1, u \in \mathbb{R}^m, v \in \mathbb{R}^n\} \quad (96)$$

*id est*, dyads may be normalized and the hull's boundary contains them;

$$\begin{aligned}\partial\{X \in \mathbb{R}^{m \times n} \mid \sum_i \sigma(X)_i \leq 1\} &= \{X \in \mathbb{R}^{m \times n} \mid \sum_i \sigma(X)_i = 1\} \\ &\supseteq \{uv^T \mid \|uv^T\| = 1, u \in \mathbb{R}^m, v \in \mathbb{R}^n\}\end{aligned}\quad (97)$$

Normalized dyads constitute the set of extreme points ([§2.6.0.0.1](#)) of this nuclear norm ball (*confer* Figure 25) which is, therefore, their convex hull.  $\square$

#### 2.3.2.0.3 Exercise. Convex hull of outer product.

Describe the interior of a Fanope.

Find the convex hull of nonorthogonal-projection matrices ([§E.1.1](#)):

$$\{UV^T \mid U \in \mathbb{R}^{N \times k}, V \in \mathbb{R}^{N \times k}, V^T U = I\} \quad (98)$$

Find the convex hull of nonsymmetric matrices bounded under some norm:

$$\{UV^T \mid U \in \mathbb{R}^{m \times k}, V \in \mathbb{R}^{n \times k}, \|UV^T\| \leq 1\} \quad (99)$$

▼

#### 2.3.2.0.4 Example. Permutation polyhedron.

[227] [354] [292]

A *permutation matrix*  $\Xi$  is formed by interchanging rows and interchanging columns of Identity matrix  $I$ . Since  $\Xi$  is square and  $\Xi^T \Xi = I$ , the set of all permutation matrices  $\Pi$  (of particular dimension) is a *proper subset* of the nonconvex *manifold* of orthogonal matrices  $\mathcal{Q}$  ([§B.5](#)). In fact, the only orthogonal matrices having all nonnegative entries are permutations of the Identity:

$$\Xi^{-1} = \Xi^T, \quad \Xi \geq \mathbf{0} \quad \Leftrightarrow \quad \Xi \in \Pi \triangleq \mathcal{Q} \cap \mathbb{R}_+ \quad (100)$$

a.k.a, the *binary orthogonal matrices*  $\Pi = \mathcal{Q} \cap \mathbb{B}$  in

$$\mathbb{B}^{n \times n} \triangleq \{0, 1\}^{n \times n} \quad (101)$$

And the only positive semidefinite permutation matrix is the Identity. [370, §6.5 prob.20]

Regarding the permutation matrices as a set of points in Euclidean space equidistant from the origin, its convex hull is a bounded polyhedron ([§2.12](#)) described (Birkhoff, 1946)

$$\begin{aligned}\mathcal{S}_\Pi \triangleq \text{conv}\{\Pi\} &= \text{conv}\{\Pi_i(I \in \mathbb{S}^n) \in \mathbb{R}^{n \times n}, i = 1 \dots n!\} \\ &= \{X \in \mathbb{R}^{n \times n} \mid X^T \mathbf{1} = \mathbf{1}, X \mathbf{1} = \mathbf{1}, X \geq \mathbf{0}\} \\ &= \{X \in \mathbb{R}^{n \times n} \mid (I \otimes \mathbf{1}^T) \text{vec } X = \mathbf{1}, (\mathbf{1}^T \otimes I) \text{vec } X = \mathbf{1}, X \geq \mathbf{0}\}\end{aligned}\quad (102)$$

where  $\Pi_i$  is a linear operator here representing the  $i^{\text{th}}$  permutation. This polyhedral hull, whose  $n!$  vertices are the permutation matrices  $\Pi$ , is known as the set of *doubly stochastic matrices* or the *permutation polyhedron*. Permutation matrices are the minimal cardinality (fewest nonzero entries) doubly stochastic matrices. The only binary matrices belonging to this polyhedron are the permutation matrices  $\Pi = \mathcal{S}_\Pi \cap \mathbb{B}$ . The only orthogonal matrices belonging to this polyhedron are permutation matrices  $\Pi = \mathcal{Q} \cap \mathcal{S}_\Pi$ .

It is remarkable that  $n!$  permutation matrices can be described as the extreme points ([§2.6.0.0.1](#)) of a bounded polyhedron, of affine dimension  $(n-1)^2$ , that is itself described by  $2n$  equalities. [2.16](#) By *Carathéodory's theorem*, conversely, any doubly stochastic matrix can be described as a convex combination of at most  $(n-1)^2 + 1$  permutation matrices. [228, §8.7] [58, thm.1.2.5] This polyhedron, then, can be a device for *relaxing* an integer, combinatorial, or Boolean optimization problem. [2.17](#) [71] [315, §3.1]  $\square$

[2.16](#)  $2n-1$  linearly independent equality constraints in  $n^2$  nonnegative variables providing  $n^2$  facets.

[2.17](#) *Relaxation* replaces an *objective function* with its *convex envelope* or expands a feasible set to one that is convex. Dantzig first showed in 1951 that, by this device, the so-called *assignment problem* can be formulated as a linear program. [353] [27, §II.5]

**2.3.2.0.5 Example.** *Convex hull of orthonormal matrices.*

[28, §1.2]

Consider rank- $k$  matrices  $U \in \mathbb{R}^{n \times k}$  such that  $U^T U = I$ . These are the orthonormal matrices; a closed bounded submanifold, of all orthogonal matrices, having dimension  $nk - \frac{1}{2}k(k+1)$  [55]. Their convex hull is expressed, for  $1 \leq k \leq n$

$$\begin{aligned} \text{conv}\{U \in \mathbb{R}^{n \times k} \mid U^T U = I\} &= \{X \in \mathbb{R}^{n \times k} \mid \|X\|_2 \leq 1\} \\ &= \{X \in \mathbb{R}^{n \times k} \mid \|X^T a\| \leq \|a\| \quad \forall a \in \mathbb{R}^n\} \end{aligned} \quad (103)$$

By Schur complement (§A.4), the *spectral norm*  $\|X\|_2$  constraining largest singular value  $\sigma(X)_1$  can be expressed as a semidefinite constraint

$$\|X\|_2 \leq 1 \Leftrightarrow \begin{bmatrix} I & X \\ X^T & I \end{bmatrix} \succeq 0 \quad (104)$$

because of equivalence  $X^T X \preceq I \Leftrightarrow \sigma(X) \preceq \mathbf{1}$  with singular values. (1719) (1604) (1605)

When  $k=n$ , matrices  $U$  are orthogonal and their convex hull is called the *spectral norm ball* which is the set of all *contractions*. [229, p.158] [365, p.313] The orthogonal matrices then constitute the extreme points (§2.6.0.0.1) of this hull. Hull intersection with the nonnegative orthant  $\mathbb{R}_+^{n \times n}$  contains the permutation polyhedron (102).  $\square$

### 2.3.3 Conic hull

In terms of a finite-length point list (or set) arranged columnar in  $X \in \mathbb{R}^{n \times N}$  (77), its conic hull is expressed

$$\mathcal{K} \triangleq \text{cone}\{x_\ell, \ell=1 \dots N\} = \text{cone } X = \{Xa \mid a \succeq 0\} \subseteq \mathbb{R}^n \quad (105)$$

*id est*, every nonnegative combination of points from the list. Conic hull of any finite-length list forms a *polyhedral cone* [225, §A.4.3] (§2.12.1.0.1); e.g, Figure 27, Figure 53a); the smallest closed convex cone (§2.7.2) that contains the list.

By convention, the aberration [366, §2.1]

$$\text{cone } \emptyset \triangleq \{\mathbf{0}\} \quad (106)$$

Given some arbitrary set  $\mathcal{C}$ , it is apparent

$$\text{conv } \mathcal{C} \subseteq \text{cone } \mathcal{C} \quad (107)$$

### 2.3.4 Vertex-description

The conditions in (79), (87), and (105) respectively define an *affine combination*, *convex combination*, and *conic combination* of elements from the set or list. Whenever a Euclidean body can be described as some hull or span of a set of points, then that representation is loosely called a *vertex-description* and those points are called *generators*.

## 2.4 Halfspace, Hyperplane

A two-dimensional affine subset is called a *plane*. An  $n-1$ -dimensional affine subset of  $\mathbb{R}^n$  is called a *hyperplane*. [343] [225] Every hyperplane partially bounds a *halfspace*.<sup>2.18</sup>

---

<sup>2.18</sup> which is convex, but not affine, and the only nonempty convex set in  $\mathbb{R}^n$  whose complement is convex and nonempty.

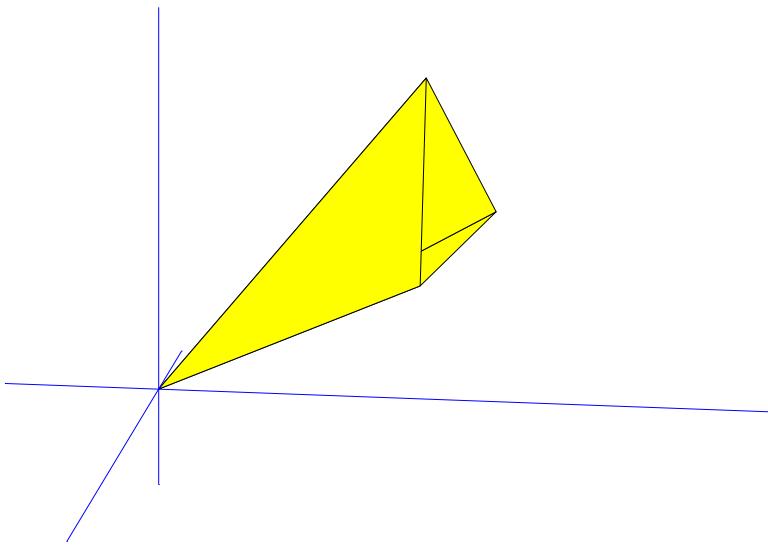


Figure 27: A simplicial cone (§2.12.3.1.1) in  $\mathbb{R}^3$  whose boundary is drawn truncated; constructed using  $A \in \mathbb{R}^{3 \times 3}$  and  $C = \mathbf{0}$  in (289). By the most fundamental definition of a cone (§2.7.1), entire boundary can be constructed from an aggregate of rays emanating exclusively from the origin. Each of three extreme directions corresponds to an edge (§2.6.0.0.3); they are conically, affinely, and linearly independent for this cone. Because this set is polyhedral, exposed directions are in one-to-one correspondence with extreme directions; there are only three. Its extreme directions give rise to what is called a *vertex-description* of this polyhedral cone; simply, the conic hull of extreme directions. Obviously this cone can also be constructed by intersection of three halfspaces; hence the equivalent *halfspace-description*.

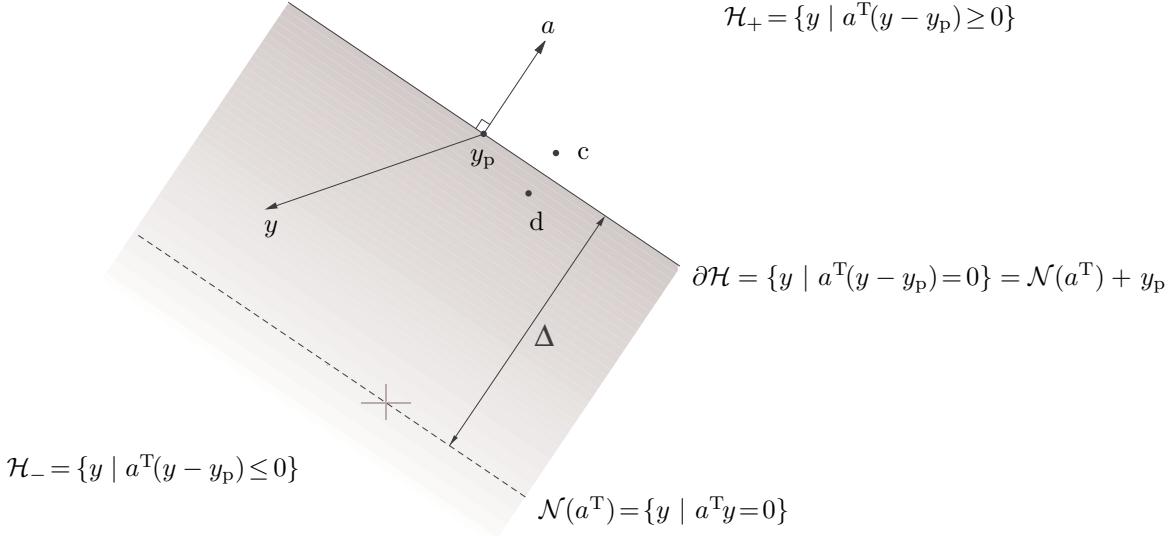


Figure 28: Hyperplane illustrated  $\partial\mathcal{H}$  is a line partially bounding halfspaces  $\mathcal{H}_-$  and  $\mathcal{H}_+$  in  $\mathbb{R}^2$ . Shaded is a rectangular piece of semiinfinite  $\mathcal{H}_-$  with respect to which vector  $a$  is outward-normal to bounding hyperplane; vector  $a$  is inward-normal with respect to  $\mathcal{H}_+$ . Halfspace  $\mathcal{H}_-$  contains nullspace  $\mathcal{N}(a^T)$  (dashed line through origin) because  $a^T y_p > 0$ . Hyperplane, halfspace, and nullspace are each drawn truncated. Points  $c$  and  $d$  are equidistant from hyperplane, and vector  $c - d$  is normal to it.  $\Delta$  is distance from origin to hyperplane.

### 2.4.1 Halfspaces $\mathcal{H}_+$ and $\mathcal{H}_-$

Euclidean space  $\mathbb{R}^n$  is partitioned in two by any hyperplane  $\partial\mathcal{H}$ ; *id est*,  $\mathcal{H}_- + \mathcal{H}_+ = \mathbb{R}^n$ . The resulting (closed convex) halfspaces, both partially bounded by  $\partial\mathcal{H}$ , may be described

$$\mathcal{H}_- = \{y \mid a^T y \leq b\} = \{y \mid a^T(y - y_p) \leq 0\} \subset \mathbb{R}^n \quad (108)$$

$$\mathcal{H}_+ = \{y \mid a^T y \geq b\} = \{y \mid a^T(y - y_p) \geq 0\} \subset \mathbb{R}^n \quad (109)$$

where nonzero vector  $a \in \mathbb{R}^n$  is an *outward-normal* to the hyperplane partially bounding  $\mathcal{H}_-$  while an *inward-normal* with respect to  $\mathcal{H}_+$ . For any vector  $y - y_p$  that makes an obtuse angle with normal  $a$ , vector  $y$  will lie in the halfspace  $\mathcal{H}_-$  on one side (shaded in Figure 28) of the hyperplane while acute angles denote  $y$  in  $\mathcal{H}_+$  on the other side.

An equivalent more intuitive representation of a halfspace comes about when we consider all the points in  $\mathbb{R}^n$  closer to point  $d$  than to point  $c$  or equidistant, in the Euclidean sense; from Figure 28,

$$\mathcal{H}_- = \{y \mid \|y - d\| \leq \|y - c\|\} \quad (110)$$

This representation, in terms of proximity, is resolved with the more conventional representation of a halfspace (108) by squaring both sides of the inequality in (110);

$$\mathcal{H}_- = \left\{ y \mid (c - d)^T y \leq \frac{\|c\|^2 - \|d\|^2}{2} \right\} = \left\{ y \mid (c - d)^T \left( y - \frac{c + d}{2} \right) \leq 0 \right\} \quad (111)$$

#### 2.4.1.1 PRINCIPLE 1: Halfspace-description of convex sets

The most fundamental principle in convex geometry follows from the *geometric Hahn-Banach theorem* [280, §5.12] [19, §1] [150, §I.1.2] which guarantees any closed convex

set to be an intersection of halfspaces.

#### 2.4.1.1.1 Theorem. Halfspaces.

[225, §A.4.2b] [43, §2.4]

A closed convex set in  $\mathbb{R}^n$  is equivalent to the intersection of all halfspaces that contain it.  $\diamond$

Intersection of multiple halfspaces in  $\mathbb{R}^n$  may be represented using a matrix constant  $A$

$$\bigcap_i \mathcal{H}_{i-} = \{y \mid A^T y \leq b\} = \{y \mid A^T(y - y_p) \leq 0\} \quad (112)$$

$$\bigcap_i \mathcal{H}_{i+} = \{y \mid A^T y \geq b\} = \{y \mid A^T(y - y_p) \geq 0\} \quad (113)$$

where  $b$  is now a vector, and the  $i^{\text{th}}$  column of  $A$  is normal to a hyperplane  $\partial\mathcal{H}_i$  partially bounding  $\mathcal{H}_i$ . By the *halfspaces theorem*, intersections like this can describe interesting convex Euclidean bodies such as polyhedra and cones (Figure 27); giving rise to the term *halfspace-description*.

## 2.4.2 Hyperplane $\partial\mathcal{H}$ representations

Every hyperplane  $\partial\mathcal{H}$  is an affine set parallel to an  $n - 1$ -dimensional subspace of  $\mathbb{R}^n$ ; it is itself a subspace if and only if it contains the origin.

$$\dim \partial\mathcal{H} = n - 1 \quad (114)$$

so a hyperplane is a point in  $\mathbb{R}$ , a line in  $\mathbb{R}^2$ , a plane in  $\mathbb{R}^3$ , and so on. Every hyperplane can be described as the intersection of complementary halfspaces; [343, §19]

$$\partial\mathcal{H} = \mathcal{H}_- \cap \mathcal{H}_+ = \{y \mid a^T y \leq b, a^T y \geq b\} = \{y \mid a^T y = b\} \quad (115)$$

a halfspace-description. Assuming *normal*  $a \in \mathbb{R}^n$  to be nonzero, then any hyperplane in  $\mathbb{R}^n$  can be described as the solution set to vector equation  $a^T y = b$  (illustrated in Figure 28 and Figure 29 for  $\mathbb{R}^2$ );

$$\partial\mathcal{H} \triangleq \{y \mid a^T y = b\} = \{y \mid a^T(y - y_p) = 0\} = \{Z\xi + y_p \mid \xi \in \mathbb{R}^{n-1}\} \subset \mathbb{R}^n \quad (116)$$

All solutions  $y$  constituting the hyperplane are *offset* from the nullspace of  $a^T$  by the same vector constant  $y_p \in \mathbb{R}^n$  that is any particular solution to  $a^T y = b$ ; *id est*,

$$y = Z\xi + y_p \quad (117)$$

where the columns of  $Z \in \mathbb{R}^{n \times n-1}$  constitute a basis for  $\mathcal{N}(a^T) = \{x \in \mathbb{R}^n \mid a^T x = 0\}$  the nullspace.<sup>2.19</sup>

Conversely, given any point  $y_p$  in  $\mathbb{R}^n$ , the unique hyperplane containing it having normal  $a$  is the affine set  $\partial\mathcal{H}$  (116) where  $b$  equals  $a^T y_p$  and where a basis for  $\mathcal{N}(a^T)$  is arranged in  $Z$  columnar. Hyperplane dimension is apparent from dimension of  $Z$ ; that hyperplane is parallel to the span of its columns.

#### 2.4.2.0.1 Exercise. Hyperplane scaling.

Given normal  $y$ , draw a hyperplane  $\{x \in \mathbb{R}^2 \mid x^T y = 1\}$ . Suppose  $z = \frac{1}{2}y$ . On the same plot, draw the hyperplane  $\{x \in \mathbb{R}^2 \mid x^T z = 1\}$ . Now suppose  $z = 2y$ , then draw the last hyperplane again with this new  $z$ . What is the apparent effect of scaling normal  $y$ ?  $\blacktriangledown$

---

<sup>2.19</sup>We will find this expression for  $y$  in terms of nullspace of  $a^T$  (more generally, of matrix  $A$  (145)) to be a useful trick (a practical device) for eliminating affine equality constraints, much as we did here.

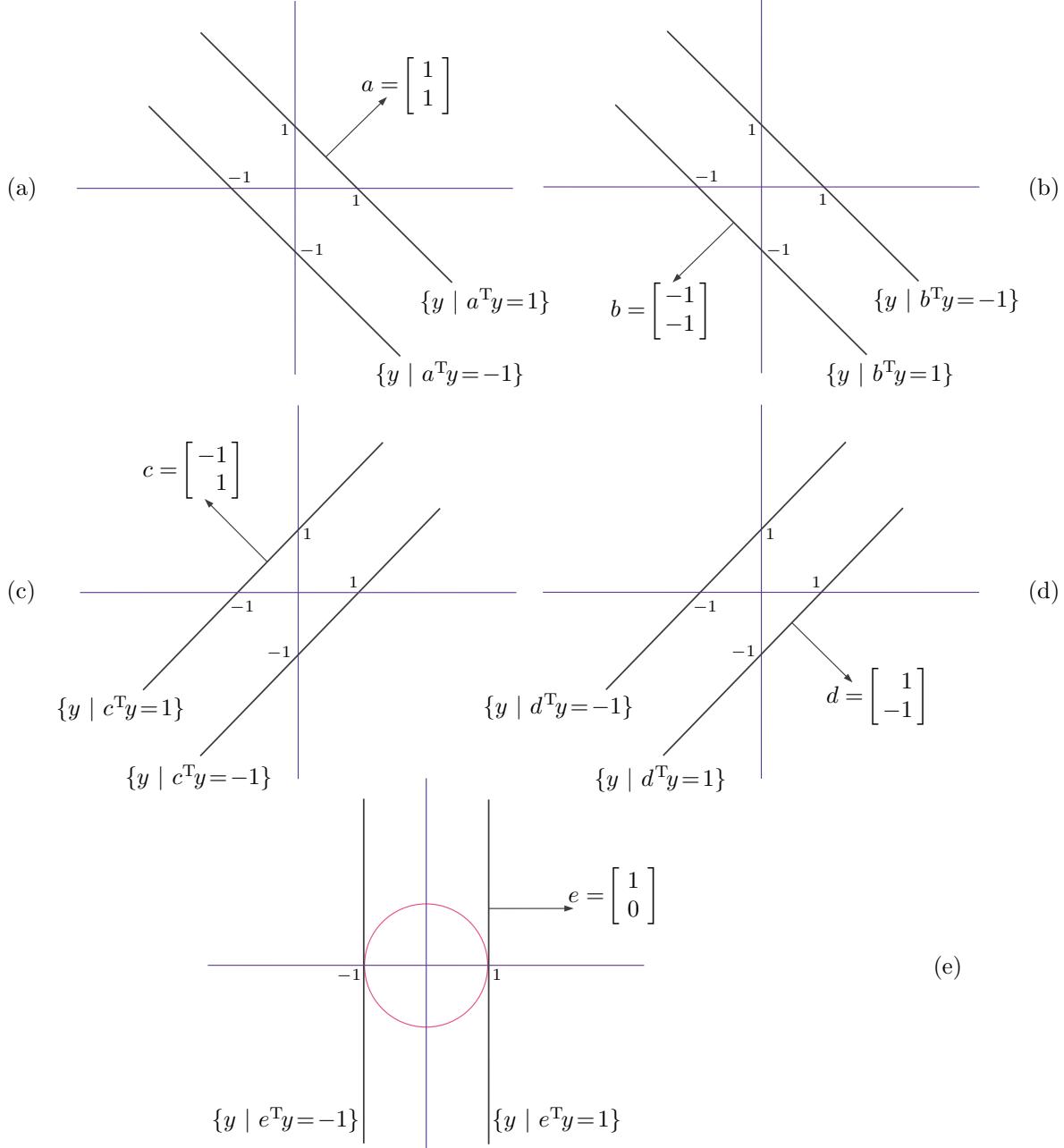


Figure 29: (a)-(d) Hyperplanes in  $\mathbb{R}^2$  (truncated) redundantly emphasize: hyperplane movement opposite to its normal direction minimizes vector inner-product. This concept is exploited to attain analytical solution of linear programs by visual inspection; e.g., §2.4.2.6.2, §2.5.1.2.2, §3.4.0.0.2, [65, exer.4.8-exer.4.20]. Each graph is also interpretable as contour plot of a real affine function of two variables as in Figure 77. (e)  $|\beta|/\|\alpha\|$  from  $\partial\mathcal{H} = \{x \mid \alpha^T x = \beta\}$  represents radius of hypersphere about  $\mathbf{0}$  supported by any hyperplane with same ratio  $|\text{inner product}|/\text{norm}$ .

#### 2.4.2.0.2 Example. Distance from origin to hyperplane.

Given the (shortest) distance  $\Delta \in \mathbb{R}_+$  from the origin to a hyperplane having normal vector  $a$ , we can find its representation  $\partial\mathcal{H}$  by dropping a perpendicular. The point thus found is the orthogonal projection of the origin on  $\partial\mathcal{H}$  (§E.5.0.0.6): equal to  $a\Delta/\|a\|$  if the origin is known *a priori* to belong to halfspace  $\mathcal{H}_-$  (Figure 28), or equal to  $-a\Delta/\|a\|$  if the origin belongs to halfspace  $\mathcal{H}_+$ ; *id est*, when  $\mathcal{H}_- \ni \mathbf{0}$

$$\partial\mathcal{H} = \{y \mid a^T(y - a\Delta/\|a\|) = 0\} = \{y \mid a^T y = \|a\|\Delta\} \quad (118)$$

or when  $\mathcal{H}_+ \ni \mathbf{0}$

$$\partial\mathcal{H} = \{y \mid a^T(y + a\Delta/\|a\|) = 0\} = \{y \mid a^T y = -\|a\|\Delta\} \quad (119)$$

Knowledge of only distance  $\Delta$  and normal  $a$  thus introduces ambiguity into the hyperplane representation.  $\square$

#### 2.4.2.1 Matrix variable

Any halfspace in  $\mathbb{R}^{mn}$  may be represented using a matrix variable. For variable  $Y \in \mathbb{R}^{m \times n}$ , given constants  $A \in \mathbb{R}^{m \times n}$  and  $b = \langle A, Y_p \rangle \in \mathbb{R}$

$$\mathcal{H}_- = \{Y \in \mathbb{R}^{mn} \mid \langle A, Y \rangle \leq b\} = \{Y \in \mathbb{R}^{mn} \mid \langle A, Y - Y_p \rangle \leq 0\} \quad (120)$$

$$\mathcal{H}_+ = \{Y \in \mathbb{R}^{mn} \mid \langle A, Y \rangle \geq b\} = \{Y \in \mathbb{R}^{mn} \mid \langle A, Y - Y_p \rangle \geq 0\} \quad (121)$$

Recall vector inner-product from §2.2:  $\langle A, Y \rangle = \text{tr}(A^T Y) = \text{vec}(A)^T \text{vec}(Y)$ .

Hyperplanes in  $\mathbb{R}^{mn}$  may, of course, also be represented using matrix variables.

$$\partial\mathcal{H} = \{Y \mid \langle A, Y \rangle = b\} = \{Y \mid \langle A, Y - Y_p \rangle = 0\} \subset \mathbb{R}^{mn} \quad (122)$$

Vector  $a$  from Figure 28 is normal to the hyperplane illustrated. Likewise, nonzero vectorized matrix  $A$  is normal to hyperplane  $\partial\mathcal{H}$ ;

$$A \perp \partial\mathcal{H} \text{ in } \mathbb{R}^{mn} \quad (123)$$

#### 2.4.2.2 Vertex-description of hyperplane

Any hyperplane in  $\mathbb{R}^n$  may be described as affine hull of a minimal set of points  $\{x_\ell \in \mathbb{R}^n, \ell = 1 \dots n\}$  arranged columnar in a matrix  $X \in \mathbb{R}^{n \times n}$ : (79)

$$\begin{aligned} \partial\mathcal{H} &= \text{aff}\{x_\ell \in \mathbb{R}^n, \ell = 1 \dots n\}, & \dim \text{aff}\{x_\ell \mid \forall \ell\} &= n-1 \\ &= \text{aff } X, & \dim \text{aff } X &= n-1 \\ &= x_1 + \mathcal{R}\{x_\ell - x_1, \ell = 2 \dots n\}, & \dim \mathcal{R}\{x_\ell - x_1, \ell = 2 \dots n\} &= n-1 \\ &= x_1 + \mathcal{R}(X - x_1 \mathbf{1}^T), & \dim \mathcal{R}(X - x_1 \mathbf{1}^T) &= n-1 \end{aligned} \quad (124)$$

where

$$\mathcal{R}(A) = \{Ax \mid \forall x\} \quad (144)$$

#### 2.4.2.3 Affine independence, minimal set

For any particular affine set, a *minimal set* of points constituting its vertex-description is an affinely independent generating set and *vice versa*.

Arbitrary given points  $\{x_i \in \mathbb{R}^n, i = 1 \dots N\}$  are *affinely independent* (a.i.) if and only if, over all  $\zeta \in \mathbb{R}^N \ni \zeta^T \mathbf{1} = 1, \zeta_k = 0 \in \mathbb{R}$  (confer §2.1.2)

$$x_i \zeta_i + \dots + x_j \zeta_j - x_k = \mathbf{0}, \quad i \neq \dots \neq j \neq k = 1 \dots N \quad (125)$$

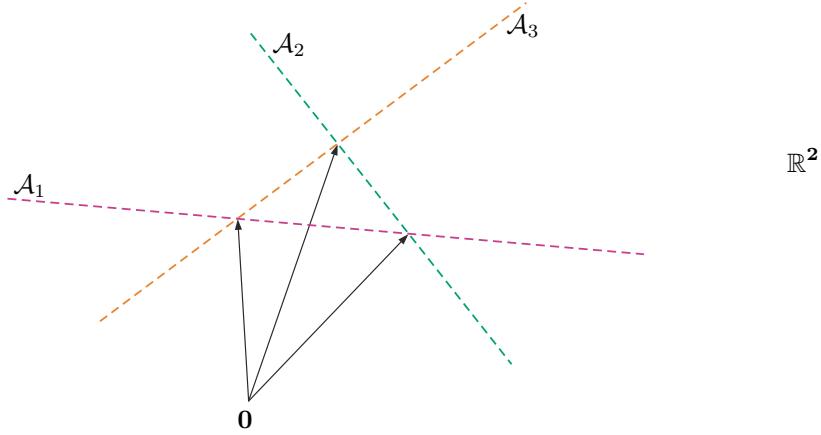


Figure 30: Of three points illustrated, any one particular point does not belong to affine hull  $\mathcal{A}_i$  ( $i \in 1, 2, 3$ , each drawn truncated) of points remaining. Three corresponding vectors are, therefore, affinely independent (but neither linearly or conically independent).

has no solution  $\zeta$ ; in words, iff no point from the given set can be expressed as an affine combination of those remaining. We deduce

$$\text{l.i.} \Rightarrow \text{a.i.} \quad (126)$$

Consequently, [232, §3] (Figure 30)

- $\{x_i, i=1 \dots N\}$  is an affinely independent set if and only if  $\{x_i - x_1, i=2 \dots N\}$  is a linearly independent (l.i.) set.

This is equivalent to the property that the columns of  $\begin{bmatrix} X \\ \mathbf{1}^T \end{bmatrix}$  (for  $X \in \mathbb{R}^{n \times N}$  as in (77)) form a linearly independent set. [225, §A.1.3]

Two nontrivial affine subsets are affinely independent iff their intersection is empty  $\{\emptyset\}$  or, analogously to subspaces, they intersect only at a point.

#### 2.4.2.4 Preservation of affine independence

Independence in the linear (§2.1.2.1), affine, and conic (§2.10.1) senses can be preserved under linear transformation. Suppose a matrix  $X \in \mathbb{R}^{n \times N}$  (77) holds an affinely independent set in its columns. Consider a transformation on the domain of such matrices

$$T(X) : \mathbb{R}^{n \times N} \rightarrow \mathbb{R}^{n \times N} \triangleq XY \quad (127)$$

where fixed matrix  $Y \triangleq [y_1 \ y_2 \ \dots \ y_N] \in \mathbb{R}^{N \times N}$  represents linear operator  $T$ . Affine independence of  $\{Xy_i \in \mathbb{R}^n, i=1 \dots N\}$  demands (by definition (125)) there exist no solution  $\zeta \in \mathbb{R}^N$  s.t.  $\zeta^T \mathbf{1} = 1$ ,  $\zeta_k = 0$ , to

$$Xy_i \zeta_i + \dots + Xy_j \zeta_j - Xy_k = \mathbf{0}, \quad i \neq \dots \neq j \neq k = 1 \dots N \quad (128)$$

By factoring out  $X$ , we see that is ensured by affine independence of  $\{y_i \in \mathbb{R}^N\}$  and by  $\mathcal{R}(Y) \cap \mathcal{N}(X) = \mathbf{0}$  where

$$\mathcal{N}(A) = \{x \mid Ax = \mathbf{0}\} \quad (145)$$

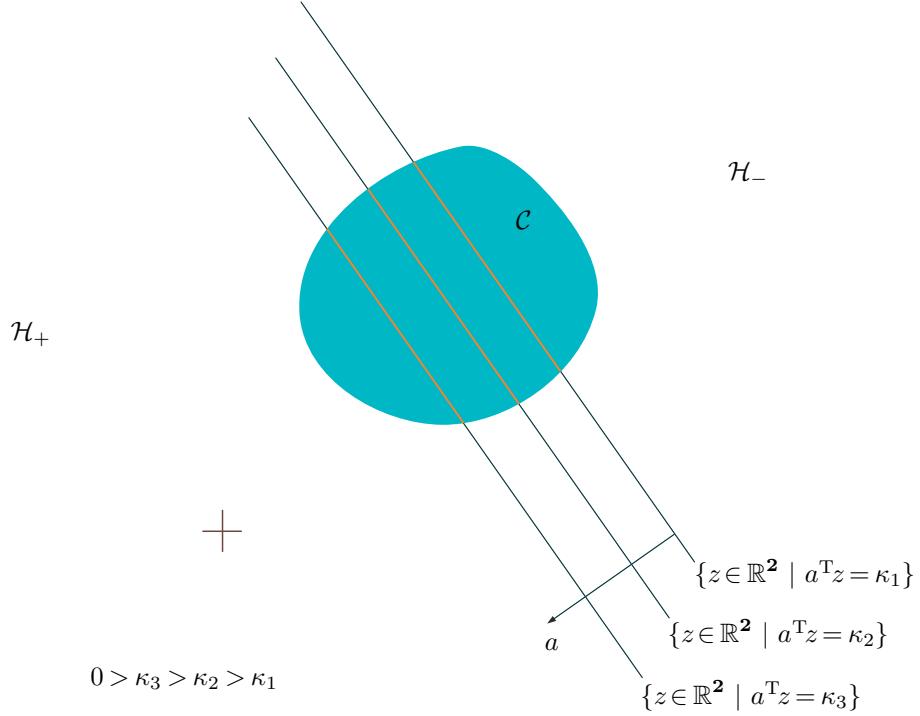


Figure 31: (confer Figure 77) Each linear contour, of equal inner product in vector  $z$  with normal  $a$ , represents  $i^{\text{th}}$  hyperplane in  $\mathbb{R}^2$  parametrized by scalar  $\kappa_i$ . Inner product  $\kappa_i$  increases in direction of normal  $a$ . In convex set  $\mathcal{C} \subset \mathbb{R}^2$ ,  $i^{\text{th}}$  line segment  $\{z \in \mathcal{C} \mid a^T z = \kappa_i\}$  represents intersection with hyperplane. (Cartesian axes for reference.)

#### 2.4.2.5 Affine maps

Affine transformations preserve affine hulls. Given any affine mapping  $T$  of vector spaces and some arbitrary set  $\mathcal{C}$  [343, p.8]

$$\text{aff}(T\mathcal{C}) = T(\text{aff } \mathcal{C}) \quad (129)$$

#### 2.4.2.6 PRINCIPLE 2: Supporting hyperplane

The second most fundamental principle of convex geometry also follows from the *geometric Hahn-Banach theorem* [280, §5.12] [19, §1] that guarantees existence of at least one hyperplane in  $\mathbb{R}^n$  supporting a full-dimensional convex set<sup>2.20</sup> at each point on its boundary.

The partial boundary  $\partial\mathcal{H}$  of a halfspace that contains arbitrary set  $\mathcal{Y}$  is called a *supporting hyperplane*  $\underline{\partial}\mathcal{H}$  to  $\mathcal{Y}$  when the hyperplane contains at least one point of  $\bar{\mathcal{Y}}$ . [343, §11]

---

<sup>2.20</sup>It is customary to speak of a hyperplane supporting set  $\mathcal{C}$  but not containing  $\mathcal{C}$ ; called *nontrivial support*. [343, p.100] Hyperplanes in support of lower-dimensional bodies are admitted.

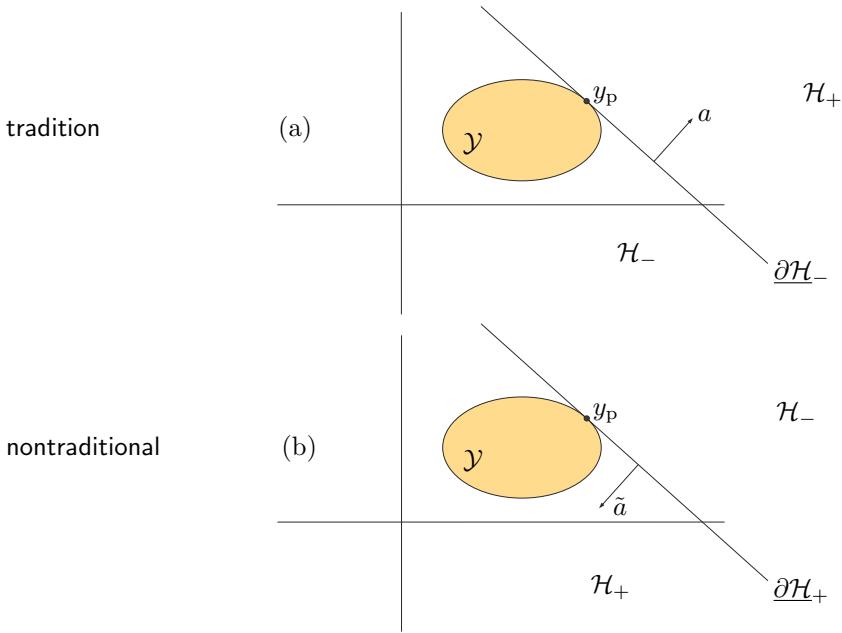


Figure 32: (a) Hyperplane  $\underline{\partial\mathcal{H}_-}$  (130) supporting closed set  $\mathcal{Y} \subset \mathbb{R}^2$ . Vector  $a$  is inward-normal to hyperplane with respect to halfspace  $\mathcal{H}_+$ , but outward-normal with respect to set  $\mathcal{Y}$ . A supporting hyperplane can be considered the limit of an increasing sequence in the normal-direction like that in Figure 31. (b) Hyperplane  $\underline{\partial\mathcal{H}_+}$  nontraditionally supporting  $\mathcal{Y}$ . Vector  $\tilde{a}$  is inward-normal to hyperplane now with respect to both halfspace  $\mathcal{H}_+$  and set  $\mathcal{Y}$ . Tradition [225] [343] recognizes only positive normal polarity in support function  $\sigma_{\mathcal{Y}}$  as in (131); *id est*, normal  $a$ , figure (a). But both interpretations of supporting hyperplane are useful.

#### 2.4.2.6.1 Definition. Supporting hyperplane $\partial\mathcal{H}$ .

Assuming set  $\mathcal{Y}$  and some normal  $a \neq \mathbf{0}$  reside in opposite halfspaces<sup>2.21</sup> (Figure 32a), then a hyperplane supporting  $\mathcal{Y}$  at point  $y_p \in \partial\mathcal{Y}$  is described

$$\underline{\partial\mathcal{H}}_- = \{y \mid a^T(y - y_p) = 0, \quad y_p \in \bar{\mathcal{Y}}, \quad a^T(z - y_p) \leq 0 \quad \forall z \in \bar{\mathcal{Y}}\} \quad (130)$$

Given only normal  $a$ , the hyperplane supporting  $\mathcal{Y}$  is equivalently described

$$\underline{\partial\mathcal{H}}_- = \{y \mid a^T y = \sup\{a^T z \mid z \in \mathcal{Y}\}\} \quad (131)$$

where real function

$$\sigma_{\mathcal{Y}}(a) = \sup\{a^T z \mid z \in \mathcal{Y}\} \quad (560)$$

is called the *support function* for  $\mathcal{Y}$ .

Another equivalent but nontraditional representation<sup>2.22</sup> for a supporting hyperplane is obtained by reversing polarity of normal  $a$ ; (1843)

$$\begin{aligned} \underline{\partial\mathcal{H}}_+ &= \{y \mid \tilde{a}^T(y - y_p) = 0, \quad y_p \in \bar{\mathcal{Y}}, \quad \tilde{a}^T(z - y_p) \geq 0 \quad \forall z \in \bar{\mathcal{Y}}\} \\ &= \{y \mid \tilde{a}^T y = -\inf\{\tilde{a}^T z \mid z \in \mathcal{Y}\} = \sup\{-\tilde{a}^T z \mid z \in \mathcal{Y}\}\} \end{aligned} \quad (132)$$

where normal  $\tilde{a}$  and set  $\mathcal{Y}$  both now reside in  $\mathcal{H}_+$  (Figure 32b).

When a supporting hyperplane contains only a single point of  $\bar{\mathcal{Y}}$ , that hyperplane is termed *strictly supporting*.<sup>2.23</sup>  $\triangle$

A full-dimensional set that has a supporting hyperplane at every point on its boundary, conversely, is convex. A convex set  $\mathcal{C} \subset \mathbb{R}^n$ , for example, can be expressed as the intersection of all halfspaces partially bounded by hyperplanes supporting it; *videlicet*, [280, p.135]

$$\bar{\mathcal{C}} = \bigcap_{a \in \mathbb{R}^n} \{y \mid a^T y \leq \sigma_{\mathcal{C}}(a)\} \quad (133)$$

by the *halfspaces theorem* (§2.4.1.1.1).

There is no geometric difference between supporting hyperplane  $\underline{\partial\mathcal{H}}_+$  or  $\underline{\partial\mathcal{H}}_-$  or  $\underline{\partial\mathcal{H}}$  and<sup>2.24</sup> an ordinary hyperplane  $\partial\mathcal{H}$  coincident with them.

#### 2.4.2.6.2 Example. Minimization over hypercube.

Consider minimization of a linear function over a *hypercube*, given vector  $c$

$$\begin{array}{ll} \text{minimize}_{x} & c^T x \\ \text{subject to} & -\mathbf{1} \preceq x \preceq \mathbf{1} \end{array} \quad (134)$$

This convex optimization problem is called a *linear program*<sup>2.25</sup> because the *objective*<sup>2.26</sup> of minimization  $c^T x$  is a linear function of variable  $x$  and the constraints describe a *polyhedron* (intersection of a finite number of halfspaces and hyperplanes).

<sup>2.21</sup>Normal  $a$  belongs to  $\mathcal{H}_+$  by definition.

<sup>2.22</sup>useful for constructing the *dual cone*; e.g., Figure 59b. Tradition would instead have us construct the *polar cone*; which is, the negative dual cone.

<sup>2.23</sup>Rockafellar terms a strictly supporting hyperplane *tangent* to  $\mathcal{Y}$  if it is unique there; [343, §18, p.169] a definition we do not adopt because our only criterion for tangency is intersection exclusively with a relative boundary. Hiriart-Urruty & Lemaréchal [225, p.44] (*confer* [343, p.100]) do not demand any tangency of a supporting hyperplane.

<sup>2.24</sup>If vector-normal polarity is unimportant, we may instead signify a supporting hyperplane by  $\partial\mathcal{H}$ .

<sup>2.25</sup>The term *program* has its roots in economics. It was originally meant with regard to a *plan* or to efficient organization or systematization of some industrial process. [103, §2]

<sup>2.26</sup>The *objective* is the function that is argument to minimization or maximization.

Any vector  $x$  satisfying the constraints is called a *feasible solution*. Applying graphical concepts from Figure 29, Figure 31, and Figure 32,  $x^* = -\text{sgn}(c)$  is an *optimal solution* to this minimization problem but is not necessarily unique. It generally holds for optimization problem solutions:

$$\text{optimal} \Rightarrow \text{feasible} \quad (135)$$

Because an optimal solution always exists at a hypercube vertex (§2.6.1.0.1) regardless of value of nonzero vector  $c$  in (134) [103, p.158] [16, p.2], mathematicians see this geometry as a means to relax a discrete problem (whose desired solution is integer or combinatorial, *confer* Example 4.2.3.1.1). [268, §3.1] [269]  $\square$

#### 2.4.2.6.3 Exercise. Unbounded below.

Suppose instead we minimize over the unit hypersphere in Example 2.4.2.6.2;  $\|x\| \leq 1$ . What is an expression for optimal solution now? Is that program still linear?

Now suppose minimization of absolute value in (134). Are the following programs equivalent for some arbitrary real convex set  $\mathcal{C}$ ? (*confer*(522))

$$\begin{array}{lll} \underset{x \in \mathbb{R}}{\text{minimize}} & |x| \\ \text{subject to} & -1 \leq x \leq 1 \\ & x \in \mathcal{C} \end{array} \equiv \begin{array}{lll} \underset{\alpha, \beta}{\text{minimize}} & \alpha + \beta \\ \text{subject to} & 1 \geq \beta \geq 0 \\ & 1 \geq \alpha \geq 0 \\ & \alpha - \beta \in \mathcal{C} \end{array} \quad (136)$$

Many optimization problems of interest and some methods of solution require nonnegative variables. The method illustrated below splits a variable into parts;  $x = \alpha - \beta$  (extensible to vectors). Under what conditions on vector  $a$  and scalar  $b$  is an optimal solution  $x^*$  negative infinity?

$$\begin{array}{lll} \underset{\alpha \in \mathbb{R}, \beta \in \mathbb{R}}{\text{minimize}} & \alpha - \beta \\ \text{subject to} & \beta \geq 0 \\ & \alpha \geq 0 \\ & a^T \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = b \end{array} \quad (137)$$

Minimization of the *objective function* entails maximization of  $\beta$ .  $\blacktriangledown$

#### 2.4.2.7 PRINCIPLE 3: Separating hyperplane

The third most fundamental principle of convex geometry again follows from the *geometric Hahn-Banach theorem* [280, §5.12] [19, §1] [150, §I.1.2] that guarantees existence of a hyperplane separating two nonempty convex sets in  $\mathbb{R}^n$  whose relative interiors are nonintersecting. *Separation* intuitively means each set belongs to a halfspace on an opposing side of the hyperplane. There are two cases of interest:

- 1) If the two sets intersect only at their relative boundaries (§2.1.7.2), then there exists a separating hyperplane  $\partial\mathcal{H}$  containing the intersection but containing no points relatively interior to either set. If at least one of the two sets is open, conversely, then the existence of a separating hyperplane implies the two sets are nonintersecting. [65, §2.5.1]
- 2) A *strictly separating hyperplane*  $\partial\mathcal{H}$  intersects the closure of neither set; its existence is guaranteed when intersection of the closures is empty and at least one set is bounded. [225, §A.4.1]

### 2.4.3 Angle between hyperspaces

Given halfspace-descriptions, *dihedral* angle between hyperplanes or halfspaces is defined as the angle between their defining normals. Given normals  $a$  and  $b$  respectively describing  $\partial\mathcal{H}_a$  and  $\partial\mathcal{H}_b$ , for example

$$\measuredangle(\partial\mathcal{H}_a, \partial\mathcal{H}_b) \triangleq \arccos\left(\frac{\langle a, b \rangle}{\|a\| \|b\|}\right) \text{ radians} \quad (138)$$

## 2.5 Subspace representations

There are two common forms of expression for Euclidean subspaces, both coming from elementary linear algebra: *range form*  $\mathcal{R}$  and *nullspace form*  $\mathcal{N}$ ; a.k.a, vertex-description and halfspace-description respectively.

The fundamental vector subspaces associated with a matrix  $A \in \mathbb{R}^{m \times n}$  [368, §3.1] are ordinarily related by orthogonal complement (Figure 18)

$$\mathcal{R}(A^T) \perp \mathcal{N}(A), \quad \mathcal{N}(A^T) \perp \mathcal{R}(A) \quad (139)$$

$$\mathcal{R}(A^T) \oplus \mathcal{N}(A) = \mathbb{R}^n, \quad \mathcal{N}(A^T) \oplus \mathcal{R}(A) = \mathbb{R}^m \quad (140)$$

and of dimension:

$$\dim \mathcal{R}(A^T) = \dim \mathcal{R}(A) = \text{rank } A \leq \min\{m, n\} \quad (141)$$

with complementarity (**a.k.a** *conservation of dimension*)

$$\dim \mathcal{N}(A) = n - \text{rank } A, \quad \dim \mathcal{N}(A^T) = m - \text{rank } A \quad (142)$$

These equations (139)-(142) comprise the *fundamental theorem of linear algebra*. [368, p.95, p.138]

From these four fundamental subspaces, the rowspace and range identify one form of subspace description (*vertex-description* (§2.3.4) or *range form*)

$$\mathcal{R}(A^T) \triangleq \text{span } A^T = \{A^T y \mid y \in \mathbb{R}^m\} = \{x \in \mathbb{R}^n \mid A^T y = x, y \in \mathcal{R}(A)\} \quad (143)$$

$$\mathcal{R}(A) \triangleq \text{span } A = \{Ax \mid x \in \mathbb{R}^n\} = \{y \in \mathbb{R}^m \mid Ax = y, x \in \mathcal{R}(A^T)\} \quad (144)$$

while the nullspaces identify the second common form (*halfspace-description* (115) or *nullspace form*)

$$\mathcal{N}(A) \triangleq \{x \in \mathbb{R}^n \mid Ax = \mathbf{0}\} = \{x \in \mathbb{R}^n \mid x \perp \mathcal{R}(A^T)\} \quad (145)$$

$$\mathcal{N}(A^T) \triangleq \{y \in \mathbb{R}^m \mid A^T y = \mathbf{0}\} = \{y \in \mathbb{R}^m \mid y \perp \mathcal{R}(A)\} \quad (146)$$

Range forms (143) (144) are realized as the respective span of the column vectors in matrices  $A^T$  and  $A$ , whereas nullspace form (145) or (146) is the solution set to a linear equation similar to hyperplane definition (116). Yet because matrix  $A$  generally has multiple rows, halfspace-description  $\mathcal{N}(A)$  is actually the intersection of as many hyperplanes through the origin; for (145), each row of  $A$  is normal to a hyperplane while each row of  $A^T$  is a normal for (146).

#### 2.5.0.0.1 Exercise. Subspace algebra.

Given

$$\mathcal{R}(A) + \mathcal{N}(A^T) = \mathcal{R}(B) + \mathcal{N}(B^T) = \mathbb{R}^m \quad (147)$$

prove

$$\mathcal{R}(A) \supseteq \mathcal{N}(B^T) \Leftrightarrow \mathcal{N}(A^T) \subseteq \mathcal{R}(B) \quad (148)$$

$$\mathcal{R}(A) \supseteq \mathcal{R}(B) \Leftrightarrow \mathcal{N}(A^T) \subseteq \mathcal{N}(B^T) \quad (149)$$

e.g, Theorem A.3.1.0.6. ▼

### 2.5.1 Subspace or affine subset...

Any particular vector subspace  $\mathcal{R}_p$  can be described as nullspace  $\mathcal{N}(A)$  of some matrix  $A$  or as range  $\mathcal{R}(B)$  of some matrix  $B$ .

More generally, we have the choice of expressing an  $n - m$ -dimensional affine subset of  $\mathbb{R}^n$  as the intersection of  $m$  hyperplanes, or as the offset span of  $n - m$  vectors:

#### 2.5.1.1 ... as hyperplane intersection

Any affine subset  $\mathcal{A}$  of dimension  $n - m$  can be described as an intersection of  $m$  hyperplanes in  $\mathbb{R}^n$ ; [343, p.6, p.44] given *wide* ( $m \leq n$ ) *full-rank* (rank =  $\min\{m, n\}$ ) matrix

$$A \triangleq \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} \in \mathbb{R}^{m \times n} \quad (150)$$

and vector  $b \in \mathbb{R}^m$ ,

$$\mathcal{A} \triangleq \{x \in \mathbb{R}^n \mid Ax = b\} = \bigcap_{i=1}^m \{x \mid a_i^T x = b_i\} \quad (151)$$

a halfspace-description. (115)

For example: The intersection of any two independent<sup>2.27</sup> hyperplanes in  $\mathbb{R}^3$  is a line, whereas three independent hyperplanes intersect at a point. Intersection of two independent hyperplanes is a plane in  $\mathbb{R}^4$  (Example 2.5.1.2.1), whereas three hyperplanes intersect at a line, four at a point, and so on.  $\mathcal{A}$  describes a subspace whenever  $b = \mathbf{0}$ .

For  $n > k$

$$\mathcal{A} \cap \mathbb{R}^k = \{x \in \mathbb{R}^n \mid Ax = b\} \cap \mathbb{R}^k = \bigcap_{i=1}^m \left\{ x \in \mathbb{R}^k \mid a_i(1:k)^T x = b_i \right\} \quad (152)$$

The result in §2.4.2.2 is extensible; *id est*, any affine subset  $\mathcal{A}$  also has a vertex-description:

#### 2.5.1.2 ... as span of nullspace basis

Alternatively, we may compute a basis for nullspace of matrix  $A$  (§E.3.1) and then equivalently express affine subset  $\mathcal{A}$  as its span plus an offset: Define

$$Z \triangleq \text{basis } \mathcal{N}(A) \in \mathbb{R}^{n \times n - \text{rank } A} \quad (153)$$

so  $AZ = \mathbf{0}$ . Then we have a vertex-description in  $Z$ ,

$$\mathcal{A} = \{x \in \mathbb{R}^n \mid Ax = b\} = \left\{ Z\xi + x_p \mid \xi \in \mathbb{R}^{n - \text{rank } A} \right\} \subseteq \mathbb{R}^n \quad (154)$$

the offset span of  $n - \text{rank } A$  column vectors, where  $x_p$  is any particular solution to  $Ax = b$ ; *e.g.*,  $\mathcal{A}$  describes a subspace whenever  $x_p = \mathbf{0}$ .

---

<sup>2.27</sup>Any number of hyperplanes are called *independent* when defining normals are linearly independent. This misuse departs from independence of two affine subsets that demands intersection only at a point or not at all. (§2.1.4.0.1)

### 2.5.1.2.1 Example. Intersecting planes in 4-space.

Two planes can intersect at a point in four-dimensional Euclidean vector space. It is easy to visualize intersection of two planes in three dimensions; a line can be formed. In four dimensions it is harder to visualize. So let's resort to the tools acquired.

Suppose an intersection of two hyperplanes in four dimensions is specified by a wide full-rank matrix  $A_1 \in \mathbb{R}^{2 \times 4}$  ( $m = 2, n = 4$ ) as in (151):

$$\mathcal{A}_1 \triangleq \left\{ x \in \mathbb{R}^4 \mid \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \end{bmatrix} x = b_1 \right\} \quad (155)$$

The nullspace of  $A_1$  is two-dimensional (from  $Z$  in (154)), so  $\mathcal{A}_1$  represents a plane in four dimensions. Similarly define a second plane in terms of  $A_2 \in \mathbb{R}^{2 \times 4}$ :

$$\mathcal{A}_2 \triangleq \left\{ x \in \mathbb{R}^4 \mid \begin{bmatrix} a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} x = b_2 \right\} \quad (156)$$

If the two planes are affinely independent and intersect, they intersect at a point because  $\begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$  is invertible;

$$\mathcal{A}_1 \cap \mathcal{A}_2 = \left\{ x \in \mathbb{R}^4 \mid \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} x = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \right\} \quad (157)$$

□

### 2.5.1.2.2 Exercise. Linear program.

Minimize a hyperplane over affine set  $\mathcal{A}$  in the nonnegative orthant

$$\begin{array}{ll} \underset{x}{\text{minimize}} & c^T x \\ \text{subject to} & Ax = b \\ & x \succeq 0 \end{array} \quad (158)$$

where  $\mathcal{A} = \{x \mid Ax = b\}$ . Two cases of interest are drawn in Figure 33. Graphically illustrate and explain optimal solutions indicated in the caption. Why is  $\alpha^*$  negative in both cases? Is there solution on the vertical axis? What causes objective unboundedness in latter case (b)? Describe all vectors  $c$  that would yield finite optimal objective in (b).

Graphical solution to linear program

$$\begin{array}{ll} \underset{x}{\text{maximize}} & c^T x \\ \text{subject to} & x \in \mathcal{P} \end{array} \quad (159)$$

is illustrated in Figure 34. Bounded set  $\mathcal{P}$  is an intersection of many halfspaces. Why is optimal solution  $x^*$  not aligned with vector  $c$  as in Cauchy-Schwarz inequality (2255)?

▼

### 2.5.1.2.3 Exercise. Nonconvex problem.

Explain why linear program

$$\begin{array}{ll} \underset{x}{\text{minimize}} & c^T x \\ \text{subject to} & Ax = b \\ & 0 \preceq x \preceq 2 \\ & x \succeq 3 \end{array} \quad (160)$$

is not convex, even though it has linear objective and affine constraints. [2.28](#)

▼

---

[2.28](#) Hint: (694).

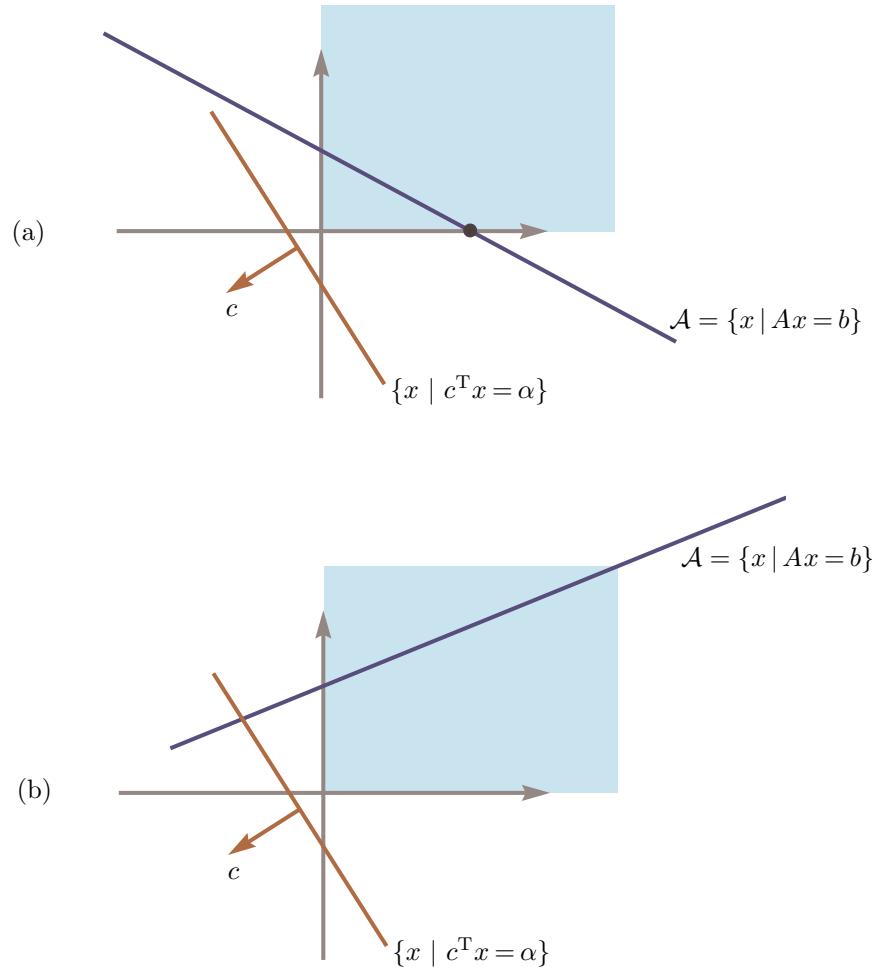


Figure 33: Minimizing hyperplane over affine subset  $\mathcal{A}$  in nonnegative orthant  $\mathbb{R}_+^2$  whose extreme directions (§2.8.1) are nonnegative Cartesian axes. Solutions visually ascertainable: (a) Optimal solution  $\bullet$  finite. (b) Optimal objective  $\alpha^* = -\infty$  unbounded.

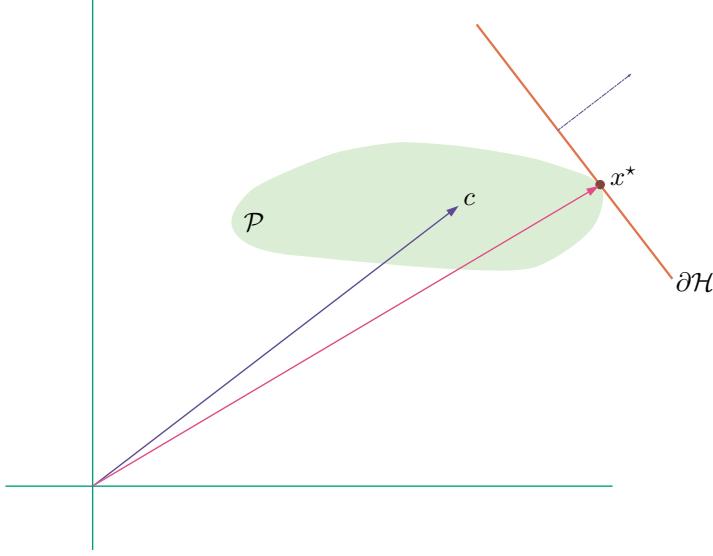


Figure 34: Maximizing hyperplane  $\partial\mathcal{H}$ , whose normal is vector  $c \in \mathcal{P}$ , over polyhedral set  $\mathcal{P}$  in  $\mathbb{R}^2$  is a linear program (159). Optimal solution  $x^*$  at  $\bullet$ .

### 2.5.2 Intersection of subspaces

The intersection of nullspaces associated with two matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{k \times n}$  can be expressed most simply as

$$\mathcal{N}(A) \cap \mathcal{N}(B) = \mathcal{N}\left(\begin{bmatrix} A \\ B \end{bmatrix}\right) \triangleq \{x \in \mathbb{R}^n \mid \begin{bmatrix} A \\ B \end{bmatrix}x = \mathbf{0}\} \quad (161)$$

nullspace of their rowwise concatenation.

Suppose the columns of a matrix  $Z$  constitute a basis for  $\mathcal{N}(A)$  while the columns of a matrix  $W$  constitute a basis for  $\mathcal{N}(BZ)$ . Then [181, §12.4.2]

$$\mathcal{N}(A) \cap \mathcal{N}(B) = \mathcal{R}(ZW) \quad (162)$$

If each basis is orthonormal, then the columns of  $ZW$  constitute an orthonormal basis for the intersection.

In the particular circumstance  $A$  and  $B$  are each positive semidefinite [22, §6], or in the circumstance  $A$  and  $B$  are two linearly independent dyads (§B.1.1), then

$$\mathcal{N}(A) \cap \mathcal{N}(B) = \mathcal{N}(A + B), \quad \begin{cases} A, B \in \mathbb{S}_+^M \\ \text{or} \\ A + B = u_1 v_1^T + u_2 v_2^T \quad (\text{l.i.}) \end{cases} \quad (163)$$

#### 2.5.2.0.1 Example. Visualization of matrix subspaces.

Fundamental subspace relations, such as

$$\mathcal{R}(A^T) \perp \mathcal{N}(A), \quad \mathcal{N}(A^T) \perp \mathcal{R}(A) \quad (139)$$

are partially defining. But to aid visualization of involved geometry, it sometimes helps to vectorize matrices; e.g., §2.9.2.5.1. For any square matrix  $A$ ,  $s \in \mathcal{N}(A)$ , and  $w \in \mathcal{N}(A^T)$

$$\langle A, ss^T \rangle = 0, \quad \langle A, ww^T \rangle = 0 \quad (164)$$

because  $s^T A s = w^T A w = 0$ . This innocuous observation can become an instrument for visualization of diagonalizable matrices (§A.5.1): For rank- $\rho$  matrix

$$A = S\Lambda S^{-1} = [s_1 \cdots s_M] \Lambda \begin{bmatrix} w_1^T \\ \vdots \\ w_M^T \end{bmatrix} = \sum_{i=1}^M \lambda_i s_i w_i^T \in \mathbb{R}^{M \times M} \quad (169)$$

nullspace eigenvectors are real (Theorem A.5.0.0.1) having range (§B.1.1)

$$\begin{aligned} \mathcal{R}\{s_i \in \mathbb{R}^M \mid \lambda_i = 0\} &= \mathcal{R}\left(\sum_{i=\rho+1}^M s_i s_i^T\right) = \mathcal{N}(A) \\ \mathcal{R}\{w_i \in \mathbb{R}^M \mid \lambda_i = 0\} &= \mathcal{R}\left(\sum_{i=\rho+1}^M w_i w_i^T\right) = \mathcal{N}(A^T) \end{aligned} \quad (165)$$

Define an unconventional basis among column vectors of each summation: (*confer*(2095))

$$\begin{aligned} \text{basis } \mathcal{N}(A) &\subseteq \sum_{i=\rho+1}^M s_i s_i^T \subseteq \mathcal{N}(A) \\ \text{basis } \mathcal{N}(A^T) &\subseteq \sum_{i=\rho+1}^M w_i w_i^T \subseteq \mathcal{N}(A^T) \end{aligned} \quad (166)$$

An *overcomplete* vectorized basis for the nullspace of any  $M \times M$  matrix is

$$\begin{aligned} \text{vec basis } \mathcal{N}(A) &= \text{vec} \sum_{i=\rho+1}^M s_i s_i^T \\ \text{vec basis } \mathcal{N}(A^T) &= \text{vec} \sum_{i=\rho+1}^M w_i w_i^T \end{aligned} \quad (167)$$

By this reckoning,  $\text{vec basis } \mathcal{R}(A) = \text{vec } A$  but is not unique. Now, because

$$\left\langle A, \sum_{i=\rho+1}^M s_i s_i^T \right\rangle = 0, \quad \left\langle A, \sum_{i=\rho+1}^M w_i w_i^T \right\rangle = 0 \quad (168)$$

then vectorized matrix  $A$  is normal to a hyperplane (of dimension  $M^2 - 1$ ) that contains both vectorized nullspaces simultaneously (each of whose dimension is  $M - \rho$ );

$$\text{vec } A \perp \text{vec basis } \mathcal{N}(A), \quad \text{vec basis } \mathcal{N}(A^T) \perp \text{vec } A \quad (169)$$

These orthogonality relations represent a departure (absent  $T$ ) from fundamental subspace relations (139) stated at the outset.  $\square$

## 2.6 Extreme, Exposed

### 2.6.0.0.1 Definition. *Extreme point.*

An extreme point  $x_\varepsilon$  of a convex set  $\mathcal{C}$  is a point, belonging to its closure  $\bar{\mathcal{C}}$  [43, §3.3], that is not expressible as a convex combination of points in  $\bar{\mathcal{C}}$  distinct from  $x_\varepsilon$ ; *id est*, for  $x_\varepsilon \in \bar{\mathcal{C}}$  and all  $x_1, x_2 \in \bar{\mathcal{C}} \setminus x_\varepsilon$

$$\mu x_1 + (1 - \mu)x_2 \neq x_\varepsilon \quad \forall \mu \in [0, 1] \quad (170)$$

$\triangle$

In other words,  $x_\varepsilon$  is an extreme point of  $\mathcal{C}$  if and only if  $x_\varepsilon$  is not a point relatively interior to any line segment in  $\bar{\mathcal{C}}$ . [399, §2.10]

Borwein & Lewis offer: [58, §4.1.6] An extreme point of a convex set  $\mathcal{C}$  is a point  $x_\varepsilon$  in  $\bar{\mathcal{C}}$  whose *relative complement*  $\bar{\mathcal{C}} \setminus x_\varepsilon$  is convex.

The set consisting of a single point  $\mathcal{C} = \{x_\varepsilon\}$  is itself an extreme point.

**2.6.0.0.2 Theorem.** *Extreme existence.* [343, §18.5.3] [27, §II.3.5]  
A nonempty closed convex set containing no lines has at least one extreme point.  $\diamond$

**2.6.0.0.3 Definition.** *Face, edge.* [225, §A.2.3]

- A *face*  $\mathcal{F}$  of convex set  $\mathcal{C}$  is a convex subset  $\mathcal{F} \subseteq \bar{\mathcal{C}}$  such that every closed line segment  $\overline{x_1x_2}$  in  $\bar{\mathcal{C}}$ , having a relatively interior point  $(x \in \text{rel intr } \overline{x_1x_2})$  in  $\mathcal{F}$ , has both endpoints in  $\mathcal{F}$ . The zero-dimensional faces of  $\mathcal{C}$  constitute its extreme points.
- All faces  $\mathcal{F}$  are extreme sets by definition; *id est*, for  $\mathcal{F} \subseteq \bar{\mathcal{C}}$  and all  $x_1, x_2 \in \bar{\mathcal{C}} \setminus \mathcal{F}$

$$\mu x_1 + (1 - \mu)x_2 \notin \mathcal{F} \quad \forall \mu \in [0, 1] \quad (171)$$

- A one-dimensional face of a convex set is called an *edge*.  $\triangle$

The empty set  $\emptyset$  and  $\bar{\mathcal{C}}$  itself are conventional faces of convex set  $\mathcal{C}$ . [343, §18] Faces of subspace  $\mathbb{R}^n$  therefore comprise only itself and  $\emptyset$ . Faces of a hyperplane  $\partial\mathcal{H}$  are constituted by itself and  $\emptyset$ . Faces of a halfspace  $\mathcal{H}$  are itself,  $\emptyset$ , and its bounding hyperplane.

Dimension of a face is the penultimate number of affinely independent points (§2.4.2.3) belonging to it;

$$\dim \mathcal{F} = \sup_{\rho} \dim \mathcal{R}\{x_2 - x_1, x_3 - x_1, \dots, x_{\rho} - x_1 \mid x_i \in \mathcal{F}, i=1 \dots \rho\} \quad (172)$$

The point of intersection in  $\bar{\mathcal{C}}$  with a strictly supporting hyperplane identifies an extreme point, but not *vice versa*. The nonempty intersection of any supporting hyperplane with  $\bar{\mathcal{C}}$  identifies a face, in general, but not *vice versa*. To acquire a converse, the concept *exposed face* requires introduction:

## 2.6.1 Exposure

**2.6.1.0.1 Definition.** *Exposed face, exposed point, vertex, facet.* [225, §A.2.3, §A.2.4]

- $\mathcal{F}$  is an *exposed face* of an  $n$ -dimensional convex set  $\mathcal{C}$  iff there is a supporting hyperplane  $\underline{\partial\mathcal{H}}$  to  $\bar{\mathcal{C}}$  such that

$$\mathcal{F} = \underline{\partial\mathcal{H}} \cap \bar{\mathcal{C}} \quad (173)$$

Only faces of dimension  $-1$  through  $n-1$  can be exposed by a hyperplane.

- An *exposed point*, the definition of *vertex*, is equivalent to a zero-dimensional exposed face; the point of intersection with a strictly supporting hyperplane.
- A *facet* is an  $n-1$ -dimensional exposed face of an  $n$ -dimensional convex set  $\mathcal{C}$ ; facets exist in one-to-one correspondence with the  $n-1$ -dimensional faces.<sup>2.29</sup>
- $\overline{\{\text{exposed points}\}} = \{\text{extreme points}\}$   
 $\overline{\{\text{exposed faces}\}} \subseteq \{\text{faces}\}$   $\triangle$

---

<sup>2.29</sup>This coincidence occurs simply because the facet's dimension is the same as dimension of the supporting hyperplane exposing it.

### 2.6.1.1 Density of exposed points

For any closed convex set  $\mathcal{C}$ , its exposed points constitute a *dense* subset of its extreme points; [343, §18] [372] [366, §3.6, p.115] dense in the sense [436] that closure of that subset yields the set of extreme points.

For the convex set illustrated in Figure 35, point B cannot be exposed because it relatively bounds both the facet  $\overline{AB}$  and the closed quarter circle, each bounding the set. Since B is not relatively interior to any line segment in the set, then B is an extreme point by definition. Point B may be regarded as the limit of some sequence of exposed points beginning at vertex C.

### 2.6.1.2 Face transitivity and algebra

Faces of a convex set enjoy transitive relation. If  $\mathcal{F}_1$  is a face (an extreme set) of  $\mathcal{F}_2$  which in turn is a face of  $\mathcal{F}_3$ , then it is always true that  $\mathcal{F}_1$  is a face of  $\mathcal{F}_3$ . (The parallel statement for exposed faces is false. [343, §18]) For example, any extreme point of  $\mathcal{F}_2$  is an extreme point of  $\mathcal{F}_3$ ; in this example,  $\mathcal{F}_2$  could be a face exposed by a hyperplane supporting polyhedron  $\mathcal{F}_3$ . [250, def.115/6 p.358] Yet it is erroneous to presume that a face, of dimension 1 or more, consists entirely of extreme points. Nor is a face of dimension 2 or more entirely composed of edges, and so on.

For the polyhedron in  $\mathbb{R}^3$  from Figure 22, for example, the nonempty faces exposed by a hyperplane are the vertices, edges, and facets; there are no more. The zero-, one-, and two-dimensional faces are in one-to-one correspondence with the exposed faces in that example.

### 2.6.1.3 Smallest face

Define the smallest face  $\mathcal{F}$ , that contains some element  $G$ , of a convex set  $\mathcal{C}$ :

$$\mathcal{F}(\mathcal{C} \ni G) \tag{174}$$

*videlicet*,  $\overline{\mathcal{C}} \supset \text{rel intr } \mathcal{F}(\mathcal{C} \ni G) \ni G$ . An affine set has no faces except itself and the empty set. The smallest face, that contains  $G$ , of intersection of convex set  $\mathcal{C}$  with an affine set  $\mathcal{A}$  [268, §2.4] [269]

$$\mathcal{F}((\mathcal{C} \cap \mathcal{A}) \ni G) = \mathcal{F}(\mathcal{C} \ni G) \cap \mathcal{A} \tag{175}$$

equals intersection of  $\mathcal{A}$  with the smallest face, that contains  $G$ , of set  $\mathcal{C}$ .

### 2.6.1.4 Conventional boundary

(confer §2.1.7.2) Relative boundary

$$\text{rel } \partial \mathcal{C} = \overline{\mathcal{C}} \setminus \text{rel intr } \mathcal{C} \tag{24}$$

is equivalent to:

#### 2.6.1.4.1 Definition. Conventional boundary of convex set. [225, §C.3.1]

The relative boundary  $\partial \mathcal{C}$  of a nonempty convex set  $\mathcal{C}$  is the union of all exposed faces of  $\overline{\mathcal{C}}$ .  $\triangle$

Equivalence to (24) comes about because it is conventionally presumed that any supporting hyperplane, central to the definition of exposure, does not contain  $\mathcal{C}$ . [343, p.100] Any face  $\mathcal{F}$  of convex set  $\mathcal{C}$  (that is not  $\mathcal{C}$  itself) belongs to  $\text{rel } \partial \mathcal{C}$ . (§2.8.2.1)

## 2.7 Cones

In optimization, convex cones achieve prominence because they generalize subspaces. Most compelling is the projection analogy: Projection on a subspace can be ascertained from projection on its orthogonal complement (Figure 192), whereas projection on a closed convex cone can be determined from projection instead on its *algebraic complement* (§2.13, Figure 193, §E.9.2); called the *polar cone*.

**2.7.0.0.1 Definition.** *Ray.*

The one-dimensional set

$$\{\zeta \Gamma + B \mid \zeta \geq 0, \Gamma \neq \mathbf{0}\} \subset \mathbb{R}^n \quad (176)$$

defines a *halfline* called a *ray* in nonzero *direction*  $\Gamma \in \mathbb{R}^n$  having *base*  $B \in \mathbb{R}^n$ . When  $B = \mathbf{0}$ , a ray is the conic hull of direction  $\Gamma$ ; hence a closed convex cone.  $\triangle$

Relative boundary of a single ray, base  $\mathbf{0}$  in any dimension, is the origin because that is the union of all exposed faces not containing the entire set. Its relative interior is the ray itself excluding the origin.

### 2.7.1 Cone defined

A set  $\mathcal{X}$  is called, simply, *cone* if and only if

$$\Gamma \in \mathcal{X} \Rightarrow \zeta \Gamma \in \overline{\mathcal{X}} \text{ for all } \zeta \geq 0 \quad (177)$$

where  $\overline{\mathcal{X}}$  denotes closure of cone  $\mathcal{X}$ ; e.g, Figure 38, Figure 39. An example of nonconvex cone is the union of two opposing quadrants:  $\mathcal{X} = \{x \in \mathbb{R}^2 \mid x_1 x_2 \geq 0\}$ . [434, §2.5] Similar examples are Figure 36 and Figure 40.

All cones obey (177) and can be defined by an aggregate of rays emanating exclusively from the origin. Hence all closed cones contain the origin  $\mathbf{0}$  but are unbounded, excepting the simplest cone  $\{\mathbf{0}\}$ . The empty set  $\emptyset$  is not a cone, but its conic hull is;

$$\text{cone } \emptyset = \{\mathbf{0}\} \quad (106)$$

### 2.7.2 Convex cone

We call set  $\mathcal{K}$  a *convex cone* iff

$$\Gamma_1, \Gamma_2 \in \mathcal{K} \Rightarrow \zeta \Gamma_1 + \xi \Gamma_2 \in \overline{\mathcal{K}} \text{ for all } \zeta, \xi \geq 0 \quad (178)$$

*id est*, if and only if any conic combination of elements from  $\mathcal{K}$  belongs to its closure. Apparent from this definition,  $\zeta \Gamma_1 \in \overline{\mathcal{K}}$  and  $\xi \Gamma_2 \in \overline{\mathcal{K}} \forall \zeta, \xi \geq 0$ ; meaning,  $\mathcal{K}$  is a cone. Set  $\mathcal{K}$  is convex since, for any particular  $\zeta, \xi \geq 0$

$$\mu \zeta \Gamma_1 + (1 - \mu) \xi \Gamma_2 \in \overline{\mathcal{K}} \quad \forall \mu \in [0, 1] \quad (179)$$

because  $\mu \zeta, (1 - \mu) \xi \geq 0$ . Obviously,

$$\{\mathcal{X}\} \supset \{\mathcal{K}\} \quad (180)$$

the set of all convex cones is a proper subset of all cones. The set of convex cones is a narrower but more familiar class of cone, any member of which can be equivalently described as the intersection of a possibly (but not necessarily) infinite number of hyperplanes (through the origin) and halfspaces whose bounding hyperplanes pass through the origin; a halfspace-description (§2.4). Convex cones need not be full-dimensional.

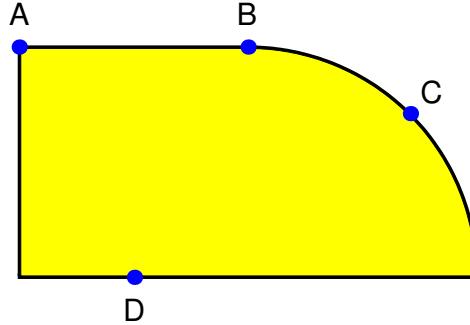


Figure 35: Closed convex set in  $\mathbb{R}^2$ . Point A is exposed hence extreme; a classical vertex. Point B is extreme but not an exposed point. Point C is exposed and extreme; zero-dimensional exposure makes it a vertex. Point D is neither an exposed or extreme point although it belongs to a one-dimensional exposed face. [225, §A.2.4] [366, §3.6] Closed face  $\overline{AB}$  is exposed; a facet. The arc is not a conventional face, yet it is composed entirely of extreme points. Union of all rotations of this entire set about its vertical edge produces another convex set in three dimensions having no edges; but that convex set produced by rotation about horizontal edge containing D has edges.

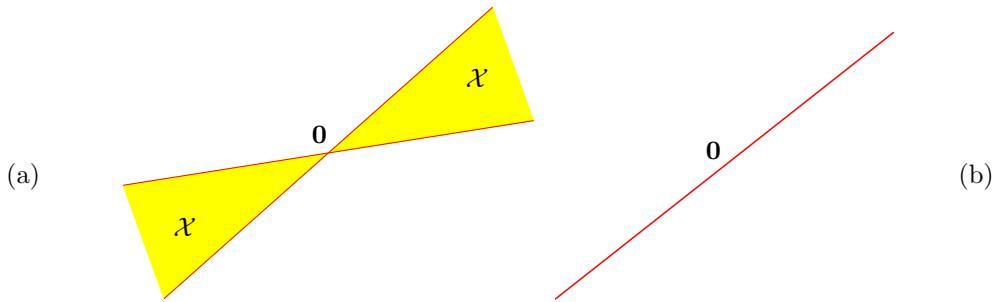


Figure 36: (a) Two-dimensional nonconvex cone drawn truncated. Boundary of this cone is itself a cone. Each half, about origin, is itself a convex cone. (b) This convex cone (drawn truncated) is a line through the origin in any dimension. It has no relative boundary, while its relative interior comprises entire line.

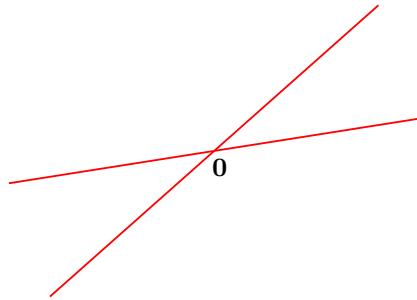


Figure 37: This nonconvex cone in  $\mathbb{R}^2$  is a pair of lines through the origin. [280, §2.4] Because the lines are linearly independent, they are algebraic complements whose vector sum is  $\mathbb{R}^2$  a convex cone.

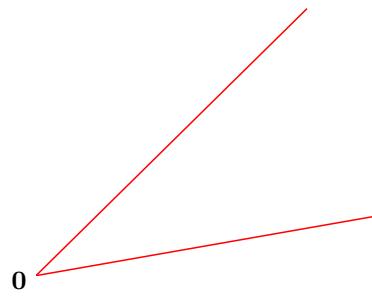


Figure 38: Boundary of a convex cone in  $\mathbb{R}^2$  is a nonconvex cone; a pair of rays emanating from the origin.

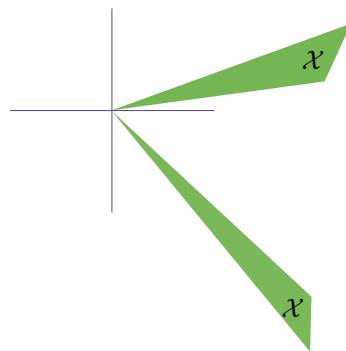


Figure 39: Union of two pointed closed convex cones in  $\mathbb{R}^2$  is nonconvex cone  $\mathcal{X}$ .

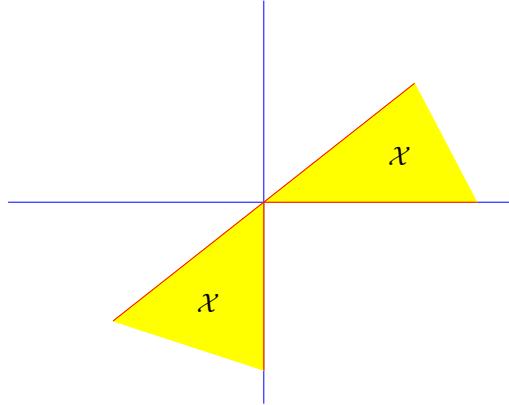


Figure 40: Truncated nonconvex cone  $\mathcal{X} = \{x \in \mathbb{R}^2 \mid x_1 \geq x_2, x_1 x_2 \geq 0\}$ . Boundary is also a cone. [280, §2.4] (Cartesian axes drawn for reference.) Each half (about the origin) is itself a convex cone.

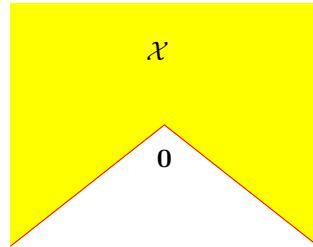


Figure 41: Nonconvex cone  $\mathcal{X}$  drawn truncated in  $\mathbb{R}^2$ . Boundary is also a cone. [280, §2.4] Cone exterior is convex cone.

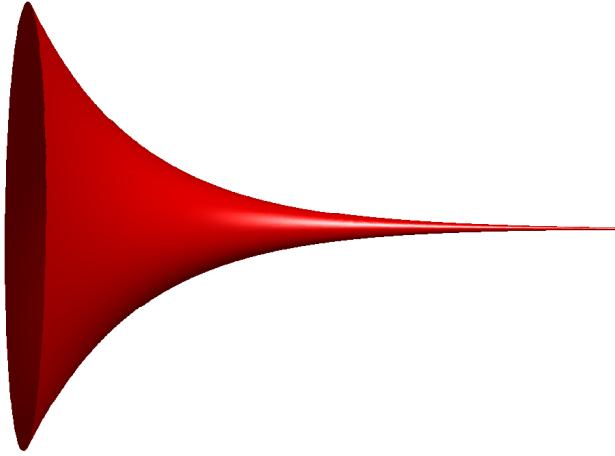


Figure 42: Not a cone; ironically, the three-dimensional *flared horn* (with or without its interior) resembling mathematical symbol  $\succ$  denoting strict cone membership and partial order.

More familiar convex cones are *Lorentz cone* (confer Figure 49)<sup>2.30</sup>

$$\mathcal{K}_\ell = \left\{ \begin{bmatrix} x \\ t \end{bmatrix} \in \mathbb{R}^n \times \mathbb{R} \mid \|x\|_\ell \leq t \right\}, \quad \ell=2 \quad (181)$$

and polyhedral cone (§2.12.1.0.1); *e.g.*, any orthant generated by Cartesian half-axes (§2.1.3). Esoteric examples of convex cones include the point at the origin, any line through the origin, any ray having the origin as base such as the nonnegative real line  $\mathbb{R}_+$  in subspace  $\mathbb{R}$ , any halfspace partially bounded by a hyperplane through the origin, the positive semidefinite cone  $\mathbb{S}_+^M$  (194), the cone of Euclidean distance matrices  $\text{EDM}^N$  (1035) (§6), *completely positive semidefinite matrices*  $\{CC^T \mid C \geq \mathbf{0}\}$  [41, p.71], any subspace, and Euclidean vector space  $\mathbb{R}^n$ .

### 2.7.2.1 cone invariance

More Euclidean bodies are cones, it seems, than are not.<sup>2.31</sup> The convex cone class of Euclidean body is invariant to scaling, linear and single- or many-valued inverse linear transformation, vector sum, and Cartesian product, but is not invariant to translation. [343, p.22]

#### 2.7.2.1.1 Theorem. Cone intersection (nonempty).

- Intersection of an arbitrary collection of convex cones is a convex cone. [343, §2, §19]
- Intersection of an arbitrary collection of closed convex cones is a closed convex cone. [289, §2.3]
- Intersection of a finite number of polyhedral cones (§2.12.1.0.1) remains a polyhedral cone.  $\diamond$

<sup>2.30</sup>a.k.a: *second-order cone*, *quadratic cone*, *circular cone* (§2.9.2.8.1), unbounded *ice-cream cone* united with its interior.

<sup>2.31</sup>confer Figures: 27 36 37 38 39 40 41 42 44 46 53 58 63 61 62 66 67 68 69 70 163 176 203

The property *pointedness* is ordinarily associated with a convex cone but, strictly speaking,

- pointed cone  $\Leftrightarrow$  convex cone (Figure 38, Figure 39)

**2.7.2.1.2 Definition.** *Pointed convex cone.* (confer §2.12.2.2)

A convex cone  $\mathcal{K}$  is *pointed* iff it contains no line. Equivalently,  $\mathcal{K}$  is not pointed iff there exists any nonzero direction  $\Gamma \in \overline{\mathcal{K}}$  such that  $-\Gamma \in \overline{\mathcal{K}}$ . If the origin is an extreme point of  $\overline{\mathcal{K}}$  or, equivalently, if

$$\overline{\mathcal{K}} \cap -\overline{\mathcal{K}} = \{\mathbf{0}\} \quad (182)$$

then  $\mathcal{K}$  is pointed, and *vice versa*. [366, §2.10] A convex cone is pointed iff the origin is the smallest nonempty face of its closure.  $\triangle$

Then a pointed closed convex cone, by principle of separating hyperplane (§2.4.2.7), has a strictly supporting hyperplane at the origin.

**2.7.2.1.3 Theorem.** *Pointed cones.* [58, §3.3.15, exer.20]

Closed convex cone  $\mathcal{K} \subset \mathbb{R}^n$  is pointed if and only if there exists a vector  $\beta$  normal to a hyperplane strictly supporting  $\mathcal{K}$ ; *id est*, for some positive scalar  $\epsilon$  (33)

$$\langle x, \beta \rangle \geq \epsilon \|x\| \quad \forall x \in \mathcal{K} \quad (183)$$

Equivalently,  $\mathcal{K}$  is pointed if and only if there exists a normal  $\alpha$  such that the set

$$\mathcal{C} \triangleq \{x \in \mathcal{K} \mid \langle x, \alpha \rangle = 1\} \quad (184)$$

is closed, bounded, and  $\mathcal{K} = \text{cone } \mathcal{C}$ .  $\diamond$

The simplest and only bounded [434, p.75] convex cone  $\mathcal{K} = \{\mathbf{0}\} \subseteq \mathbb{R}^n$  is pointed, by convention, but generally not full-dimensional. Its relative boundary is the empty set  $\emptyset$  (25) while its relative interior is the point  $\mathbf{0}$  itself (12). The pointed closed convex cone that is a halfline, emanating from the origin in  $\mathbb{R}^n$ , has relative boundary  $\mathbf{0}$  while its relative interior is the halfline itself excluding  $\mathbf{0}$ .

Pointed are any Lorentz cone, cone of Euclidean distance matrices  $\mathbb{EDM}^N$  in symmetric hollow subspace  $\mathbb{S}_h^N$ , and positive semidefinite cone  $\mathbb{S}_+^M$  in ambient  $\mathbb{S}^M$ .

If closed convex cone  $\mathcal{K}$  is not pointed, then it has no extreme point.<sup>2.32</sup> Yet a pointed closed convex cone has only one extreme point [43, §3.3]: the exposed point residing at the origin; its vertex. Pointedness is invariant to Cartesian product by (182). And from the *cone intersection theorem* it follows that an intersection of convex cones is pointed if at least one of the cones is; implying, each and every nonempty exposed face of a pointed closed convex cone is a pointed closed convex cone.

---

<sup>2.32</sup> nor does it have extreme directions (§2.8.1).

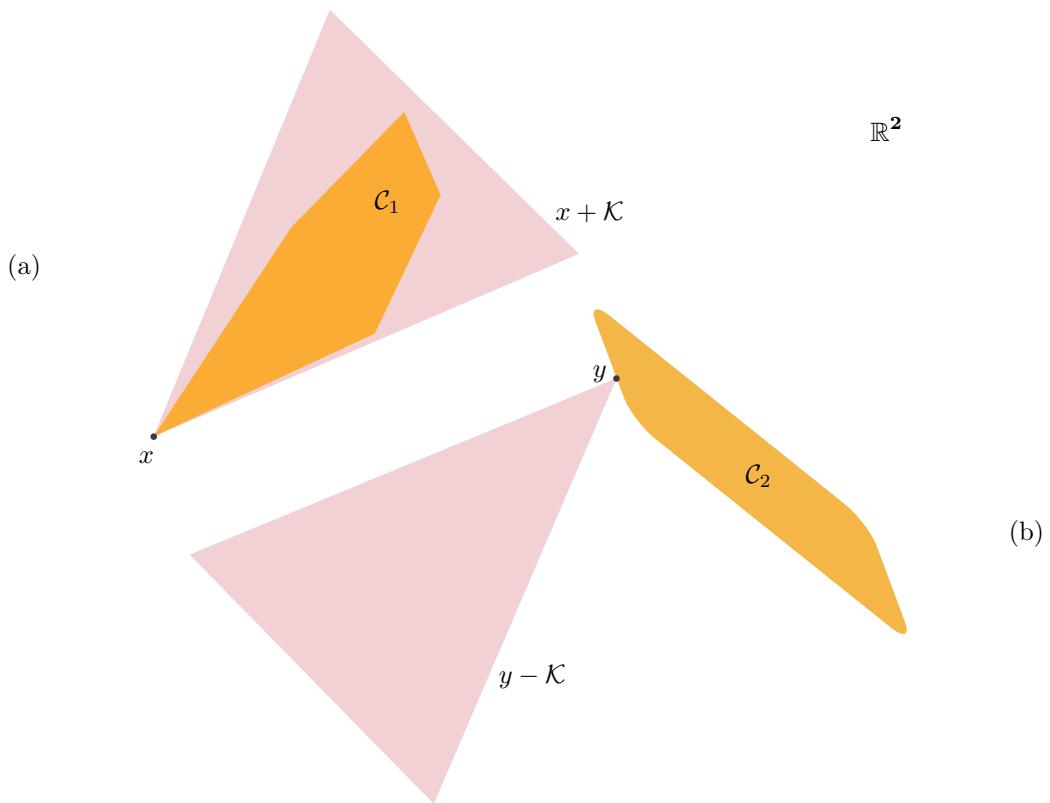


Figure 43: (confer Figure 73) (a) Point  $x$  is the unique minimum element of set  $\mathcal{C}_1$  with respect to pointed closed convex cone  $\mathcal{K}$  because cone, translated to  $x \in \mathcal{C}_1$ , contains entire set. (Cones drawn truncated.) (b) Point  $y$  is a minimal element of set  $\mathcal{C}_2$  with respect to cone  $\mathcal{K}$  because negative cone translated to  $y \in \mathcal{C}_2$  contains only  $y$ . These two concepts, *minimum/minimal*, become equivalent under a *total order*.

### 2.7.2.2 Pointed closed convex cone induces partial order

Relation  $\preceq$  represents *partial order* on some set if that relation possesses<sup>2.33</sup>

reflexivity  $(x \preceq x)$

antisymmetry  $(x \preceq z, z \preceq x \Rightarrow x = z)$

transitivity  $(x \preceq y, y \preceq z \Rightarrow x \preceq z)$ ,

$(x \preceq y, y \prec z \Rightarrow x \prec z)$

A pointed closed convex cone  $\mathcal{K}$  induces partial order on  $\mathbb{R}^n$  or  $\mathbb{R}^{m \times n}$ , [22, §1] [360, p.7] essentially defined by vector or matrix inequality;

$$x \underset{\mathcal{K}}{\preceq} z \Leftrightarrow z - x \in \mathcal{K} \quad (185)$$

$$x \underset{\mathcal{K}}{\prec} z \Leftrightarrow z - x \in \text{rel intr } \mathcal{K} \quad (186)$$

Neither  $x$  or  $z$  is necessarily a member of  $\mathcal{K}$  for these relations to hold. Only when  $\mathcal{K}$  is a nonnegative orthant  $\mathbb{R}_+^n$  do these inequalities reduce to ordinary entrywise comparison (§2.13.4.2.3) while partial order lingers. Inclusive of that special case, we ascribe nomenclature *generalized inequality* to comparison with respect to a pointed closed convex cone.

We say two points  $x$  and  $y$  are *comparable* when  $x \preceq y$  or  $y \preceq x$  with respect to pointed closed convex cone  $\mathcal{K}$ . Visceral mechanics of actually comparing points, when cone  $\mathcal{K}$  is not an orthant, are well illustrated in the example of Figure 67 which relies on the equivalent membership-interpretation in definition (185) or (186).

Comparable points and the *minimum element* of some vector- or matrix-valued partially ordered set are thus well defined, so nonincreasing sequences with respect to cone  $\mathcal{K}$  can therefore converge in this sense: Point  $x \in \mathcal{C}$  is the unique minimum element of set  $\mathcal{C}$  with respect to cone  $\mathcal{K}$  iff for each and every  $z \in \mathcal{C}$  we have  $x \preceq z$ ; equivalently, iff  $\mathcal{C} \subseteq x + \mathcal{K}$ .<sup>2.34</sup>

A closely related concept, *minimal element*, is useful for partially ordered sets having no minimum element: Point  $x \in \mathcal{C}$  is a minimal element of set  $\mathcal{C}$  with respect to pointed closed convex cone  $\mathcal{K}$  if and only if  $(x - \mathcal{K}) \cap \mathcal{C} = x$ . (Figure 43) No uniqueness is implied here, although implicit is the assumption:  $\dim \mathcal{K} \geq \dim \text{aff } \mathcal{C}$ . In words, a point that is a minimal element is smaller (with respect to  $\mathcal{K}$ ) than any other point in the set to which it is comparable.

Further properties of partial order with respect to pointed closed convex cone  $\mathcal{K}$  are not defining:

homogeneity  $(x \preceq y, \lambda \geq 0 \Rightarrow \lambda x \preceq \lambda z)$ ,  $(x \prec y, \lambda > 0 \Rightarrow \lambda x \prec \lambda z)$

additivity  $(x \preceq z, u \preceq v \Rightarrow x + u \preceq z + v)$ ,  $(x \prec z, u \preceq v \Rightarrow x + u \prec z + v)$

**2.33** A set is *totally ordered* if it further obeys a comparability property of the relation: for each and every  $x$  and  $y$  from the set,  $x \preceq y$  or  $y \preceq x$ ; e.g., one-dimensional real vector space  $\mathbb{R}$  is the smallest unbounded totally ordered and connected set.

**2.34** Borwein & Lewis [58, §3.3 exer.21] ignore possibility of equality to  $x + \mathcal{K}$  in this condition, and require a second condition: ... and  $\mathcal{C} \subset y + \mathcal{K}$  for some  $y$  in  $\mathbb{R}^n$  implies  $x \in y + \mathcal{K}$ .

**2.7.2.2.1 Definition.** *Proper cone:* a cone that is

- pointed
- closed
- convex
- full-dimensional.

△

A proper cone remains proper under injective linear transformation. [98, §5.1] Examples of proper cones are the positive semidefinite cone  $\mathbb{S}_+^M$  in the ambient space of symmetric matrices (§2.9), the nonnegative real line  $\mathbb{R}_+$  in vector space  $\mathbb{R}$ , or any orthant in  $\mathbb{R}^n$ , and the set of all coefficients of univariate degree- $n$  polynomials nonnegative on interval  $[0, 1]$  [65, exmp.2.16] or univariate degree- $2n$  polynomials nonnegative over  $\mathbb{R}$  [65, exer.2.37].

## 2.8 Cone boundary

Every hyperplane supporting a convex cone contains the origin. [225, §A.4.2] Because any supporting hyperplane to a convex cone must therefore itself be a cone, then from the *cone intersection theorem* (§2.7.2.1.1) it follows:

**2.8.0.0.1 Lemma.** *Cone faces.*

[27, §II.8]

Each nonempty exposed face of a convex cone is a convex cone. ◇

**2.8.0.0.2 Theorem.** *Proper-cone boundary.*

Suppose a nonzero point  $\Gamma$  lies on the boundary  $\partial\mathcal{K}$  of proper cone  $\mathcal{K}$  in  $\mathbb{R}^n$ . Then it follows that the ray  $\{\zeta\Gamma \mid \zeta \geq 0\}$  also belongs to  $\partial\mathcal{K}$ . ◇

**Proof.** By virtue of its propriety, a proper cone guarantees existence of a strictly supporting hyperplane at the origin. [343, cor.11.7.3]<sup>2.35</sup> Hence the origin belongs to the boundary of  $\mathcal{K}$  because it is the zero-dimensional exposed face. The origin belongs to the ray through  $\Gamma$ , and the ray belongs to  $\mathcal{K}$  by definition (177). By the *cone faces lemma*, each and every nonempty exposed face must include the origin. Hence the closed line segment  $\overline{0\Gamma}$  must lie in an exposed face of  $\mathcal{K}$  because both endpoints do by Definition 2.6.1.4.1. That means there exists a supporting hyperplane  $\underline{\partial\mathcal{H}}$  to  $\mathcal{K}$  containing  $\overline{0\Gamma}$ . So the ray through  $\Gamma$  belongs both to  $\mathcal{K}$  and to  $\underline{\partial\mathcal{H}}$ .  $\underline{\partial\mathcal{H}}$  must therefore expose a face of  $\mathcal{K}$  that contains the ray; *id est*,

$$\{\zeta\Gamma \mid \zeta \geq 0\} \subseteq \mathcal{K} \cap \underline{\partial\mathcal{H}} \subset \partial\mathcal{K} \quad (187)$$

◆

Proper cone  $\{\mathbf{0}\}$  in  $\mathbb{R}^0$  has no boundary (24) because (12)

$$\text{rel intr}\{\mathbf{0}\} = \{\mathbf{0}\} \quad (188)$$

The boundary of any proper cone in  $\mathbb{R}$  is the origin.

The boundary of any convex cone whose dimension exceeds 1 can be constructed entirely from an aggregate of rays emanating exclusively from the origin.

---

<sup>2.35</sup>Rockafellar's corollary yields a supporting hyperplane at the origin to any convex cone in  $\mathbb{R}^n$  not equal to  $\mathbb{R}^n$ .

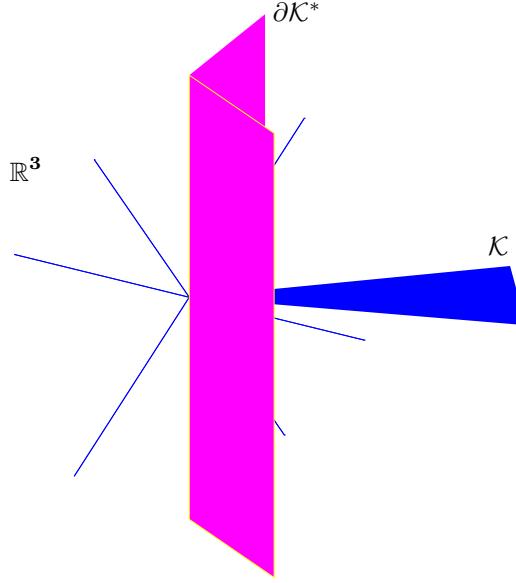


Figure 44:  $\mathcal{K}$  is a pointed polyhedral cone but not full-dimensional (drawn truncated in a plane parallel to the ground). Dual cone  $\mathcal{K}^*$  is a *wedge* having no extreme direction and no vertex; a nonpointed polyhedral cone whose truncated boundary is illustrated (faces drawn perpendicular to ground). In this particular instance,  $\mathcal{K} \subset \text{intr } \mathcal{K}^*$  excepting the origin. (Cartesian coordinate axes drawn for reference.)

### 2.8.1 Extreme direction

The property *extreme direction* arises naturally in connection with the pointed closed convex cone  $\mathcal{K} \subset \mathbb{R}^n$ , being analogous to extreme point. [343, §18, p.162]<sup>2.36</sup> An extreme direction  $\Gamma_\varepsilon$  of pointed  $\mathcal{K}$  is a vector corresponding to an edge that is a ray  $\{\zeta \Gamma_\varepsilon \in \mathcal{K} \mid \zeta \geq 0\}$  emanating from the origin.<sup>2.37</sup> Nonzero direction  $\Gamma_\varepsilon$  in pointed  $\mathcal{K}$  is extreme if and only if

$$\zeta_1 \Gamma_1 + \zeta_2 \Gamma_2 \neq \Gamma_\varepsilon \quad \forall \zeta_1, \zeta_2 \geq 0, \quad \forall \Gamma_1, \Gamma_2 \in \mathcal{K} \setminus \{\zeta \Gamma_\varepsilon \in \mathcal{K} \mid \zeta \geq 0\} \quad (189)$$

In words, an extreme direction in a pointed closed convex cone is the direction of a ray (called an *extreme ray*) that cannot be expressed as a conic combination of directions of any rays in the cone distinct from it.

An extreme ray is a one-dimensional face of  $\mathcal{K}$ . By (107), extreme direction  $\Gamma_\varepsilon$  is not a point relatively interior to any line segment in  $\mathcal{K} \setminus \{\zeta \Gamma_\varepsilon \in \mathcal{K} \mid \zeta \geq 0\}$ . Thus, by analogy, the corresponding extreme ray  $\{\zeta \Gamma_\varepsilon \in \mathcal{K} \mid \zeta \geq 0\}$  is not a ray relatively interior to any *plane segment*<sup>2.38</sup> in  $\mathcal{K}$ .

#### 2.8.1.1 extreme distinction, uniqueness

An extreme *direction* is unique, but its vector representation  $\Gamma_\varepsilon$  is not because any positive scaling of it produces another vector in the same (extreme) direction. Hence an extreme direction is *unique* to within a positive scaling. When we say extreme directions are

<sup>2.36</sup>We diverge from Rockafellar's extreme direction: "extreme point at infinity".

<sup>2.37</sup>An edge (§2.6.0.0.3) of a convex cone is not necessarily a ray. A convex cone may contain an edge that is a line; *e.g.*, a wedge-shaped polyhedral cone ( $\mathcal{K}^*$  in Figure 44).

<sup>2.38</sup>A planar fragment; in this context, a planar cone.

*distinct*, we are referring to distinctness of rays containing them. Nonzero vectors of various length in the same extreme direction are therefore interpreted to be identical extreme directions.<sup>2.39</sup>

The extreme directions of the polyhedral cone in Figure 27 (p.59), for example, correspond to its three edges. For any pointed polyhedral cone, there is a one-to-one correspondence of one-dimensional faces with extreme directions.

Extreme directions of the positive semidefinite cone (§2.9) comprise the infinite set of all symmetric rank-one matrices. [22, §6] [221, §III] It is sometimes prudent to instead consider the less infinite but complete normalized set of extreme directions: for  $M > 0$  (confer(240))

$$\{zz^T \in \mathbb{S}^M \mid \|z\|=1\} \quad (190)$$

The positive semidefinite cone in one dimension  $M=1$ ,  $\mathbb{S}_+$  the nonnegative real line, has one extreme direction belonging to its relative interior; an idiosyncrasy of dimension 1.

Pointed convex cone  $\mathcal{K}=\{\mathbf{0}\}$  has an extreme point but no extreme direction because extreme directions are nonzero by definition.

- If closed convex cone  $\mathcal{K}$  is not pointed, then its nonempty faces comprise no extreme directions and no vertex. [22, §1]

Conversely, pointed closed convex cone  $\mathcal{K}$  is equivalent to the convex hull of its vertex and all its extreme directions. [343, §18, p.167] That is the practical utility of extreme direction; to facilitate construction of polyhedral sets, apparent from the *extremes theorem*:

**2.8.1.1.1 Theorem.** (Klee) *Extremes.* [366, §3.6] [343, §18, p.166]  
(confer §2.3.2, §2.12.2.0.1) Any closed convex set containing no lines can be expressed as the convex hull of its extreme points and extreme rays. ◇

It follows that any element of a convex set containing no lines may be expressed as a linear combination of its extreme elements; e.g., §2.9.2.7.1.

### 2.8.1.2 generators

In the narrowest sense, *generators* for a convex set comprise any collection of points and directions whose convex hull constructs the set.

When the *extremes theorem* applies, the extreme points and directions are called generators of a convex set. An arbitrary collection of generators for a convex set includes its extreme elements as a subset; the set of extreme elements of a convex set is a minimal set of generators for that convex set. Any polyhedral set has a minimal set of generators whose *cardinality* is finite.

When the convex set under scrutiny is a closed convex cone, conic combination of generators during construction is implicit as shown in Example 2.8.1.2.1 and Example 2.10.2.0.1. So, a vertex at the origin (if it exists) becomes benign.

We can, of course, generate affine sets by taking the affine hull of any collection of points and directions. We broaden, thereby, the meaning of generator to be inclusive of all kinds of hulls.

Any hull of generators is loosely called a *vertex-description*. (§2.3.4) Hulls encompass subspaces, so any basis constitutes generators for a vertex-description; span basis  $\mathcal{R}(A)$ .

---

<sup>2.39</sup>Like vectors, an extreme direction can be identified with the Cartesian point at the vector's head with respect to the origin.

### 2.8.1.2.1 Example. Application of extremes theorem.

Given an extreme point at the origin and  $N$  extreme rays  $\{\zeta \Gamma_i, i=1\dots N \mid \zeta \geq 0\}$  (§2.7.0.0.1), denoting the  $i^{\text{th}}$  extreme direction by  $\Gamma_i \in \mathbb{R}^n$ , then their convex hull (87) is

$$\begin{aligned}\mathcal{P} &= \{[\mathbf{0} \ \Gamma_1 \ \Gamma_2 \dots \ \Gamma_N] a \zeta \mid a^T \mathbf{1} = 1, a \succeq 0, \zeta \geq 0\} \\ &= \{[\Gamma_1 \ \Gamma_2 \dots \ \Gamma_N] a \zeta \mid a^T \mathbf{1} \leq 1, a \succeq 0, \zeta \geq 0\} \\ &= \{[\Gamma_1 \ \Gamma_2 \dots \ \Gamma_N] b \mid b \succeq 0\} \subset \mathbb{R}^n\end{aligned}\quad (191)$$

a closed convex set that is simply a conic hull like (105).  $\square$

## 2.8.2 Exposed direction

**2.8.2.0.1 Definition.** Exposed point & direction of pointed convex cone. [343, §18] (confer §2.6.1.0.1)

- When a convex cone has a vertex, an exposed point, it resides at the origin; there can be only one.
- In the closure of a pointed convex cone, an *exposed direction* is the direction of a one-dimensional exposed face that is a ray emanating from the origin.
- $\{\text{exposed directions}\} \subseteq \{\text{extreme directions}\}$   $\triangle$

For a proper cone in vector space  $\mathbb{R}^n$  with  $n \geq 2$ , we can say more:

$$\overline{\{\text{exposed directions}\}} = \{\text{extreme directions}\} \quad (192)$$

It follows from Lemma 2.8.0.0.1 for any pointed closed convex cone, there is one-to-one correspondence of one-dimensional exposed faces with exposed directions; *id est*, there is no one-dimensional exposed face that is not a ray base  $\mathbf{0}$ .

The pointed closed convex cone  $\mathbb{EDM}^2$ , for example, is a ray in isomorphic subspace  $\mathbb{R}$  whose relative boundary (§2.6.1.4.1) is the origin. The conventionally exposed directions of  $\mathbb{EDM}^2$  constitute the empty set  $\emptyset \subset \{\text{extreme direction}\}$ . This cone has one extreme direction belonging to its relative interior; an idiosyncrasy of dimension 1.

### 2.8.2.1 Connection between boundary and extremes

**2.8.2.1.1 Theorem.** Exposed. [343, §18.7] (confer §2.8.1.1.1)  
Any closed convex set  $\mathcal{C}$  containing no lines (and whose dimension is at least 2) can be expressed as closure of the convex hull of its exposed points and exposed rays.  $\diamond$

From Theorem 2.8.1.1.1,

$$\begin{aligned}\text{rel } \partial \mathcal{C} &= \overline{\mathcal{C}} \setminus \text{rel intr } \mathcal{C} \\ &= \overline{\text{conv}\{\text{exposed points and exposed rays}\}} \setminus \text{rel intr } \mathcal{C} \\ &= \overline{\text{conv}\{\text{extreme points and extreme rays}\}} \setminus \text{rel intr } \mathcal{C}\end{aligned}\quad \left. \begin{array}{l} (24) \\ \end{array} \right\} \quad (193)$$

Thus each and every extreme point of a convex set (that is not a point) resides on its relative boundary, while each and every extreme direction of a convex set (that is not a halfline and contains no line) resides on its relative boundary because extreme points and directions of such respective sets do not belong to relative interior by definition.

The relationship between extreme sets and the relative boundary actually goes deeper: Any face  $\mathcal{F}$  of convex set  $\mathcal{C}$  (that is not  $\mathcal{C}$  itself) belongs to  $\text{rel } \partial \mathcal{C}$ , so  $\dim \mathcal{F} < \dim \mathcal{C}$ . [343, §18.1.3]

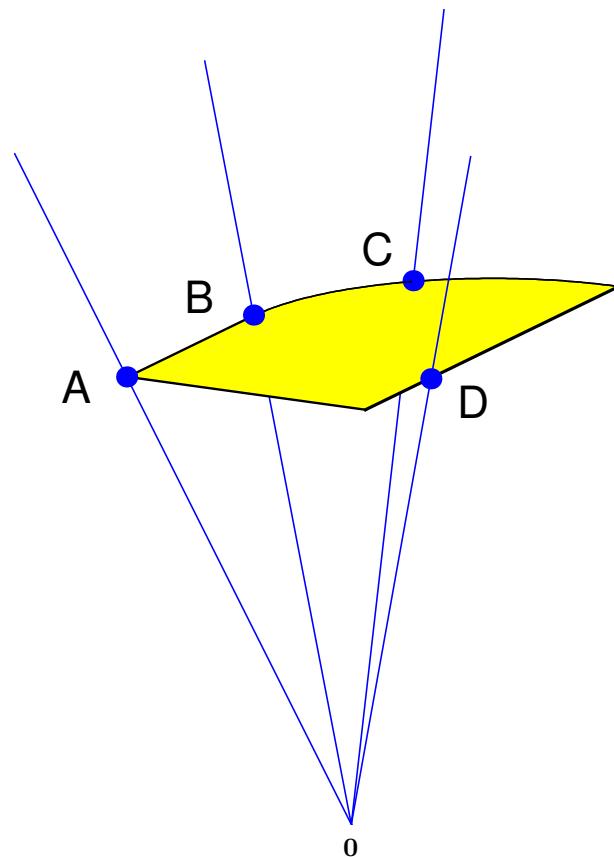


Figure 45: Properties of extreme points carry over to extreme directions. [343, §18] Four rays (drawn truncated) on boundary of conic hull of two-dimensional closed convex set from Figure 35 lifted to  $\mathbb{R}^3$ . Ray through point A is exposed hence extreme. Extreme direction B on cone boundary is not an exposed direction, although it belongs to the exposed face  $\text{cone}\{A, B\}$ . Extreme ray through C is exposed. Point D is neither an exposed or extreme direction although it belongs to a two-dimensional exposed face of the conic hull.

### 2.8.2.2 Converse *caveat*

It is inconsequent to presume that each and every extreme point and direction is necessarily exposed, as might be erroneously inferred from the *conventional boundary definition* (§2.6.1.4.1); although it can correctly be inferred: each and every extreme point and direction belongs to some exposed face.

Arbitrary points residing on the relative boundary of a convex set are not necessarily exposed or extreme points. Similarly, the direction of an arbitrary ray, base  $\mathbf{0}$ , on the boundary of a convex cone is not necessarily an exposed or extreme direction. For the polyhedral cone illustrated in Figure 27, for example, there are three two-dimensional exposed faces constituting the entire boundary, each composed of an infinity of rays. Yet there are only three exposed directions.

Neither is an extreme direction on the boundary of a pointed convex cone necessarily an exposed direction. Lift the two-dimensional set in Figure 35, for example, into three dimensions such that no two points in the set are collinear with the origin. Then its conic hull can have an extreme direction  $B$  on the boundary that is not an exposed direction, illustrated in Figure 45.

## 2.9 Positive semidefinite (PSD) cone

*The cone of positive semidefinite matrices studied in this section is arguably the most important of all non-polyhedral cones whose facial structure we completely understand.*

— Alexander Barvinok [27, p.78]

### 2.9.0.0.1 Definition. Positive semidefinite cone.

The set of all symmetric positive semidefinite matrices of particular dimension  $M$  is called the *positive semidefinite cone*:

$$\begin{aligned} \mathbb{S}_+^M &\triangleq \left\{ A \in \mathbb{S}^M \mid A \succeq 0 \right\} \\ &= \left\{ A \in \mathbb{S}^M \mid y^T A y \geq 0 \quad \forall \|y\| = 1 \right\} \\ &= \bigcap_{\|y\|=1} \left\{ A \in \mathbb{S}^M \mid \langle yy^T, A \rangle \geq 0 \right\} \\ &\equiv \{A \in \mathbb{S}_+^M \mid \text{rank } A \leq M\} \end{aligned} \tag{194}$$

formed by the intersection of an infinite number of halfspaces (§2.4.1.1) in vectorized variable<sup>2.40</sup>  $A$ , each halfspace having partial boundary containing the origin in isomorphic  $\mathbb{R}^{M(M+1)/2}$ . It is a unique immutable proper cone (§2.7.2.2.1) in the ambient space of symmetric matrices  $\mathbb{S}^M$ .

The positive definite (full-rank) matrices comprise the cone interior

$$\begin{aligned} \text{intr } \mathbb{S}_+^M &= \left\{ A \in \mathbb{S}^M \mid A \succ 0 \right\} \\ &= \left\{ A \in \mathbb{S}^M \mid y^T A y > 0 \quad \forall \|y\| = 1 \right\} \\ &= \bigcap_{\|y\|=1} \left\{ A \in \mathbb{S}^M \mid \langle yy^T, A \rangle > 0 \right\} \\ &= \{A \in \mathbb{S}_+^M \mid \text{rank } A = M\} \end{aligned} \tag{195}$$

<sup>2.40</sup> infinite in number when  $M > 1$ . Because  $y^T A y = y^T A^T y$ , matrix  $A$  is almost always assumed symmetric. (§A.2.1)

while all singular positive semidefinite matrices (having at least one 0 eigenvalue) reside on the cone boundary (Figure 46); (§A.7.5)

$$\begin{aligned}\partial \mathbb{S}_+^M &= \left\{ A \in \mathbb{S}^M \mid A \succeq 0, A \not\succ 0 \right\} \\ &= \left\{ A \in \mathbb{S}^M \mid \min\{\lambda(A)_i, i=1 \dots M\} = 0 \right\} \\ &= \left\{ A \in \mathbb{S}_+^M \mid \langle yy^T, A \rangle = 0 \text{ for some } \|y\| = 1 \right\} \\ &= \{A \in \mathbb{S}_+^M \mid \text{rank } A < M\}\end{aligned}\tag{196}$$

where  $\lambda(A) \in \mathbb{R}^M$  holds the eigenvalues of  $A$ .  $\triangle$

The only symmetric positive semidefinite matrix in  $\mathbb{S}_+^M$  having  $M$  0-eigenvalues resides at the origin. (§A.7.3.0.1)

### 2.9.0.1 Membership

Observe notation  $A \succeq 0$  denoting a positive semidefinite matrix;<sup>2.41</sup> meaning (*confer* §2.3.1.1), matrix  $A$  belongs to the positive semidefinite cone in the subspace of symmetric matrices whereas  $A \succ 0$  denotes membership to that cone's interior. (§2.13.2) Notation  $A \succ 0$ , denoting a positive definite matrix, can be read: *symmetric matrix  $A$  exceeds the origin with respect to the positive semidefinite cone interior*. These notations further imply that coordinates [*sic*] for orthogonal expansion of a positive (semi)definite matrix must be its (nonnegative) positive eigenvalues (§2.13.8.1.1, §E.6.4.1.1) when expanded in its *eigenmatrices* (§A.5.0.3); *id est*, eigenvalues must be (nonnegative) positive.

Generalizing comparison on the real line, the notation  $A \succeq B$  denotes comparison with respect to the positive semidefinite cone; (§A.3.1) *id est*,

$$A \succeq B \Leftrightarrow A - B \in \mathbb{S}_+^M\tag{197}$$

but neither matrix  $A$  or  $B$  necessarily belongs to the positive semidefinite cone. Yet, (1640)  $A \succeq B, B \succeq 0 \Rightarrow A \succeq 0$ ; *id est*,  $A \in \mathbb{S}_+^M$ . (*confer* Figure 67)

#### 2.9.0.1.1 Example. Equality constraints in semidefinite program (697).

Employing properties of partial order (§2.7.2.2) for the pointed closed convex positive semidefinite cone, it is easy to show, given  $A + S = C$

$$\begin{aligned}S \succeq 0 &\Leftrightarrow A \preceq C \\ S \succ 0 &\Leftrightarrow A \prec C\end{aligned}\tag{198}$$

$\square$

## 2.9.1 Positive semidefinite cone is convex

The set of all positive semidefinite matrices forms a convex cone in the ambient space of symmetric matrices because any pair satisfies definition (178); [228, §7.1] *videlicet*, for all  $\zeta_1, \zeta_2 \geq 0$  and each and every  $A_1, A_2 \in \mathbb{S}^M$

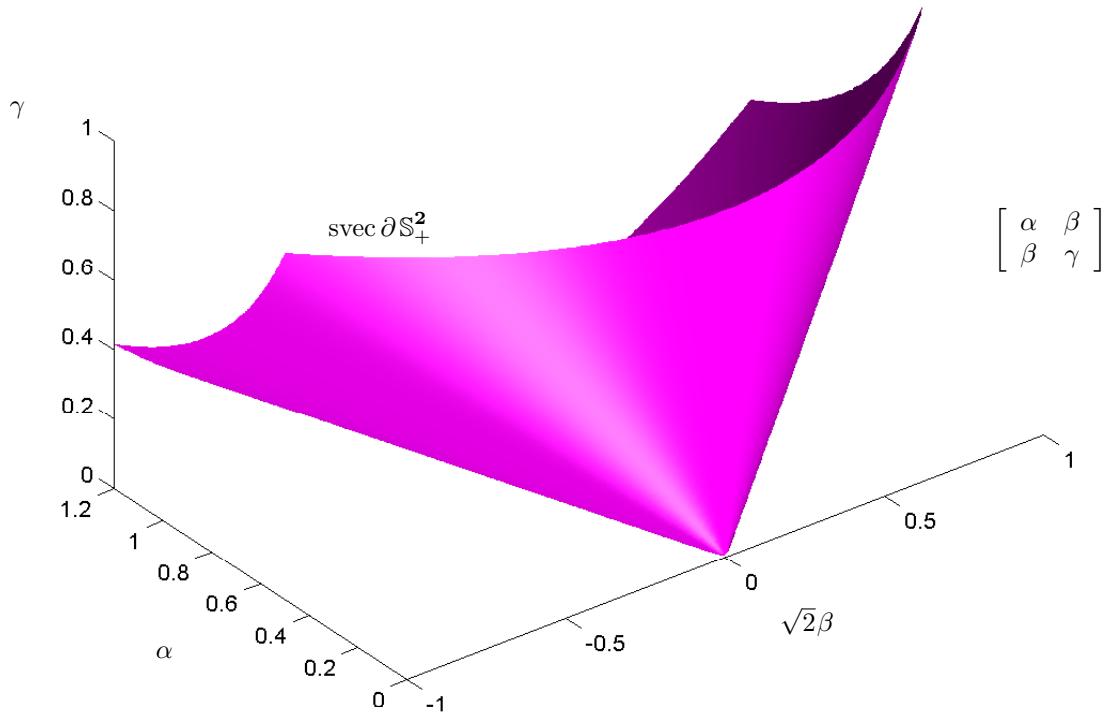
$$\zeta_1 A_1 + \zeta_2 A_2 \succeq 0 \Leftrightarrow A_1 \succeq 0, A_2 \succeq 0\tag{199}$$

a fact easily verified by the definitive test for positive semidefiniteness of a symmetric matrix (§A):

$$A \succeq 0 \Leftrightarrow x^T A x \geq 0 \text{ for each and every } \|x\| = 1\tag{200}$$

---

<sup>2.41</sup> the same as *nonnegative definite matrix*.



Minimal set of generators are the extreme directions:  $\text{svec}\{yy^T \mid y \in \mathbb{R}^M\}$

Figure 46: (d'Aspremont) Truncated boundary of PSD cone in  $\mathbb{S}^2$  plotted in isometrically isomorphic  $\mathbb{R}^3$  via svec (57); 0-contour of smallest eigenvalue (196). Lightest shading is closest, darkest shading is farthest and inside shell. Entire boundary can be constructed from an aggregate of rays (§2.7.0.0.1) emanating exclusively from origin:  $\{\kappa^2[z_1^2 \sqrt{2}z_1 z_2 z_2^2]^T \mid \kappa \in \mathbb{R}, z \in \mathbb{R}^2\}$ . A circular cone in this dimension (§2.9.2.8), each and every ray on boundary corresponds to an extreme direction but such is not the case in any higher dimension (*confer* Figure 27). PSD cone geometry is not as simple in higher dimensions [27, §II.12] although PSD cone is selfdual (383) in ambient real space of symmetric matrices. [221, §II] PSD cone has no two-dimensional face in any dimension, its only extreme point residing at  $\mathbf{0}$ .

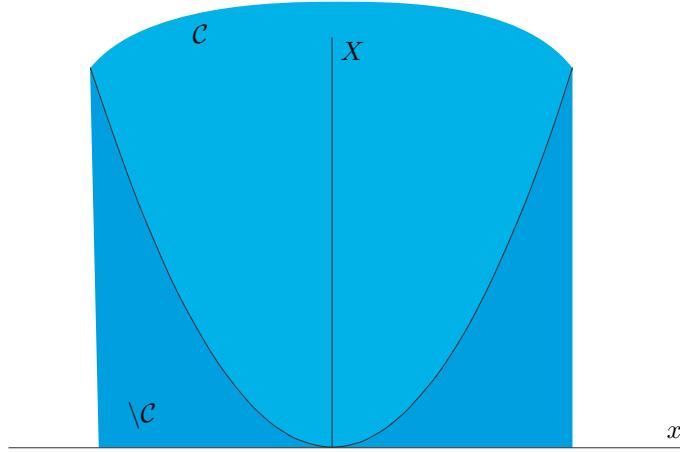


Figure 47: Convex set  $\mathcal{C} = \{X \in \mathbb{S} \times x \in \mathbb{R} \mid X \succeq xx^T\}$  drawn truncated.

*id est*, for  $A_1, A_2 \succeq 0$  and each and every  $\zeta_1, \zeta_2 \geq 0$

$$\zeta_1 x^T A_1 x + \zeta_2 x^T A_2 x \geq 0 \quad \text{for each and every normalized } x \in \mathbb{R}^M \quad (201)$$

The convex cone  $\mathbb{S}_+^M$  is more easily visualized in the isomorphic vector space  $\mathbb{R}^{M(M+1)/2}$  whose dimension is the number of free variables in a symmetric  $M \times M$  matrix. When  $M=2$  the PSD cone is semiinfinite in expanse in  $\mathbb{R}^3$ , having boundary illustrated in Figure 46. When  $M=3$  the PSD cone is six-dimensional, and so on.

#### 2.9.1.0.1 Example. Sets from maps of positive semidefinite cone.

The set

$$\mathcal{C} = \{X \in \mathbb{S}^n \times x \in \mathbb{R}^n \mid X \succeq xx^T\} \quad (202)$$

is convex because it has *Schur-form*; (§A.4)

$$X - xx^T \succeq 0 \Leftrightarrow f(X, x) \triangleq \begin{bmatrix} X & x \\ x^T & 1 \end{bmatrix} \succeq 0 \quad (203)$$

e.g., Figure 47. Set  $\mathcal{C}$  is the inverse image (§2.1.9.0.1) of  $\mathbb{S}_+^{n+1}$  under affine mapping  $f$ . The set  $\setminus \mathcal{C} = \{X \in \mathbb{S}^n \times x \in \mathbb{R}^n \mid X \prec xx^T\}$  is not convex, in contrast, having no Schur-form. Yet for fixed  $x = x_p$ , the set

$$\{X \in \mathbb{S}^n \mid X \preceq x_p x_p^T\} \quad (204)$$

is simply the negative semidefinite cone shifted to  $x_p x_p^T$ .  $\square$

#### 2.9.1.0.2 Example. Inverse image of positive semidefinite cone.

Now consider finding the set of all matrices  $X \in \mathbb{S}^N$  satisfying

$$AX + B \succeq 0 \quad (205)$$

given  $A, B \in \mathbb{S}^N$ . Define the set

$$\mathcal{X} \triangleq \{X \mid AX + B \succeq 0\} \subseteq \mathbb{S}^N \quad (206)$$

which is the inverse image of the positive semidefinite cone under affine transformation  $g(X) \triangleq AX + B$ . Set  $\mathcal{X}$  must therefore be convex by Theorem 2.1.9.0.1.

Yet we would like a less amorphous characterization of this set, so instead we consider its vectorization (37) which is easier to visualize:

$$\text{vec } g(X) = \text{vec}(AX) + \text{vec } B = (I \otimes A) \text{vec } X + \text{vec } B \quad (207)$$

where

$$I \otimes A \triangleq Q\Lambda Q^T \in \mathbb{S}^{N^2} \quad (208)$$

is block-diagonal formed by *Kronecker product* (§A.1.1 no.33, §D.1.2.1). Assign

$$\begin{aligned} x &\triangleq \text{vec } X \in \mathbb{R}^{N^2} \\ b &\triangleq \text{vec } B \in \mathbb{R}^{N^2} \end{aligned} \quad (209)$$

then make the equivalent problem: Find

$$\text{vec } \mathcal{X} = \{x \in \mathbb{R}^{N^2} \mid (I \otimes A)x + b \in \mathcal{K}\} \quad (210)$$

where

$$\mathcal{K} \triangleq \text{vec } \mathbb{S}_+^N \quad (211)$$

is a proper cone isometrically isomorphic with the positive semidefinite cone in the subspace of symmetric matrices; the vectorization of every element of  $\mathbb{S}_+^N$ . Utilizing the diagonalization (208),

$$\begin{aligned} \text{vec } \mathcal{X} &= \{x \mid \Lambda Q^T x \in Q^T(\mathcal{K} - b)\} \\ &= \{x \mid \Phi Q^T x \in \Lambda^\dagger Q^T(\mathcal{K} - b)\} \subseteq \mathbb{R}^{N^2} \end{aligned} \quad (212)$$

where  $^\dagger$  denotes matrix *pseudoinverse* (§E) and

$$\Phi \triangleq \Lambda^\dagger \Lambda \quad (213)$$

is a diagonal projection matrix whose entries are either 1 or 0 (§E.3). We have the complementary sum

$$\Phi Q^T x + (I - \Phi)Q^T x = Q^T x \quad (214)$$

So, adding  $(I - \Phi)Q^T x$  to both sides of the membership within (212) admits

$$\begin{aligned} \text{vec } \mathcal{X} &= \{x \in \mathbb{R}^{N^2} \mid Q^T x \in \Lambda^\dagger Q^T(\mathcal{K} - b) + (I - \Phi)Q^T x\} \\ &= \{x \mid Q^T x \in \Phi(\Lambda^\dagger Q^T(\mathcal{K} - b)) \oplus (I - \Phi)\mathbb{R}^{N^2}\} \\ &= \{x \in Q\Lambda^\dagger Q^T(\mathcal{K} - b) \oplus Q(I - \Phi)\mathbb{R}^{N^2}\} \\ &= (I \otimes A)^\dagger(\mathcal{K} - b) \oplus \mathcal{N}(I \otimes A) \end{aligned} \quad (215)$$

where we used the facts: linear function  $Q^T x$  in  $x$  on  $\mathbb{R}^{N^2}$  is a bijection, and  $\Phi\Lambda^\dagger = \Lambda^\dagger$ .

$$\text{vec } \mathcal{X} = (I \otimes A)^\dagger \text{vec}(\mathbb{S}_+^N - B) \oplus \mathcal{N}(I \otimes A) \quad (216)$$

In words, set  $\text{vec } \mathcal{X}$  is the vector sum of the translated PSD cone (linearly mapped onto the rowspace of  $I \otimes A$  (§E)) and the nullspace of  $I \otimes A$  (synthesis of fact from §A.6.3 and §A.7.3.0.1). Should  $I \otimes A$  have no nullspace, then  $\text{vec } \mathcal{X} = (I \otimes A)^{-1} \text{vec}(\mathbb{S}_+^N - B)$  which is the expected result.  $\square$

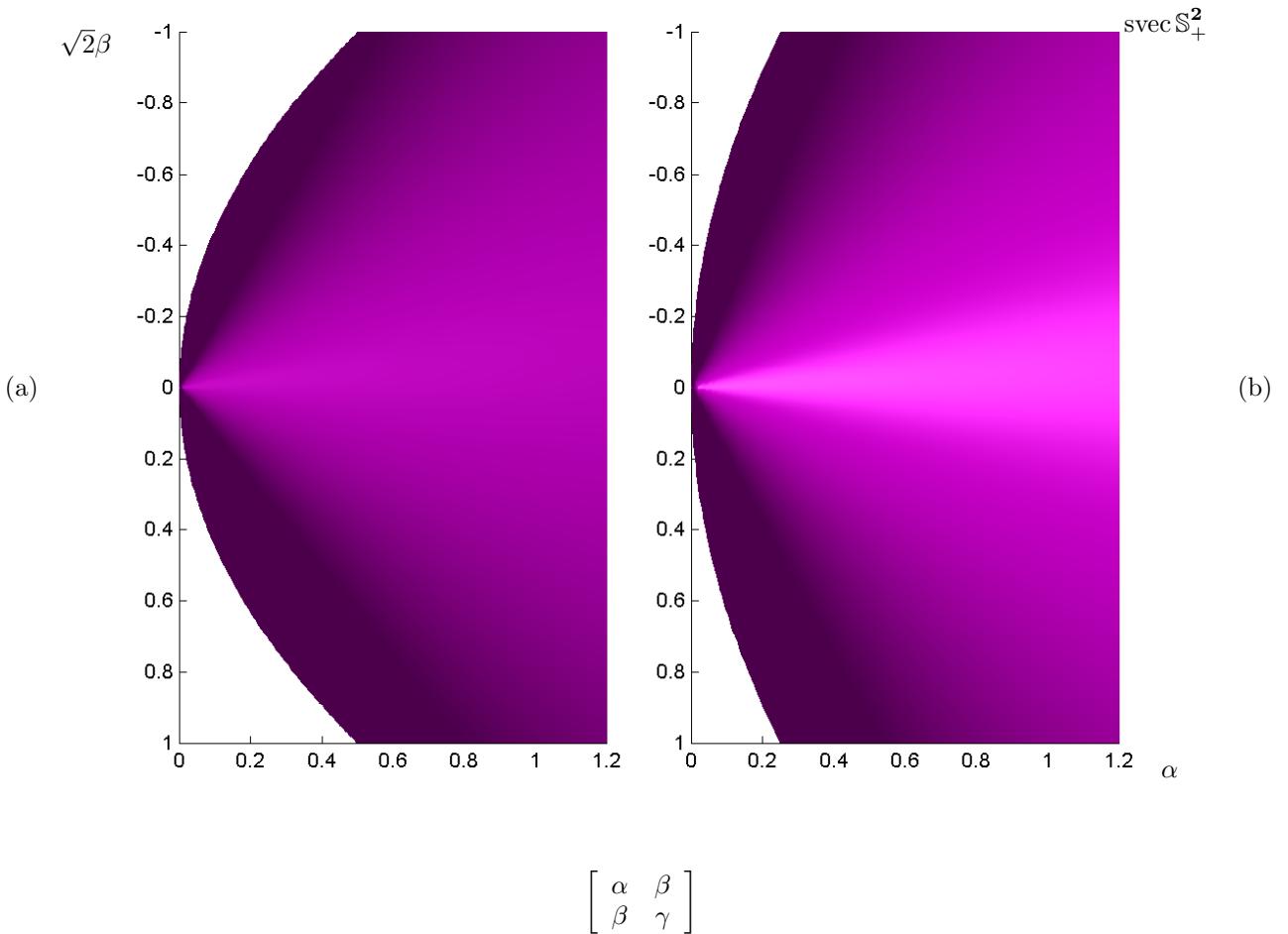


Figure 48: **(a)** Projection of truncated PSD cone  $\mathbb{S}_+^2$ , truncated above  $\gamma=1$ , on  $\alpha\beta$ -plane in isometrically isomorphic  $\mathbb{R}^3$ . View is from above with respect to Figure 46. **(b)** Truncated above  $\gamma=2$ . From these plots we might infer, for example, line  $\{[0 \ 1/\sqrt{2} \ \gamma]^T \mid \gamma \in \mathbb{R}\}$  intercepts PSD cone at some large value of  $\gamma$ ; in fact,  $\gamma=\infty$ .

### 2.9.2 Positive semidefinite cone boundary

For any symmetric positive semidefinite matrix  $A$  of rank  $\rho$ , there must exist a rank  $\rho$  matrix  $Y$  such that  $A$  be expressible as an outer product in  $Y$ ; [368, §6.3]

$$A = YY^T \in \mathbb{S}_+^M, \quad \text{rank } A = \text{rank } Y = \rho, \quad Y \in \mathbb{R}^{M \times \rho} \quad (217)$$

Then the boundary of the positive semidefinite cone may be expressed

$$\partial \mathbb{S}_+^M = \left\{ A \in \mathbb{S}_+^M \mid \text{rank } A < M \right\} = \left\{ YY^T \mid Y \in \mathbb{R}^{M \times M-1} \right\} \quad (218)$$

Because the boundary of any convex body is obtained with closure of its relative interior (§2.1.7, §2.1.7.2), from (195) we must also have

$$\begin{aligned} \mathbb{S}_+^M &= \overline{\left\{ A \in \mathbb{S}_+^M \mid \text{rank } A = M \right\}} = \overline{\left\{ YY^T \mid Y \in \mathbb{R}^{M \times M}, \text{rank } Y = M \right\}} \\ &= \overline{\left\{ YY^T \mid Y \in \mathbb{R}^{M \times M} \right\}} \end{aligned} \quad (219)$$

#### 2.9.2.1 rank $\rho$ subset of the positive semidefinite cone

For the same reason (closure), this applies more generally; for  $0 \leq \rho \leq M$

$$\overline{\left\{ A \in \mathbb{S}_+^M \mid \text{rank } A = \rho \right\}} = \left\{ A \in \mathbb{S}_+^M \mid \text{rank } A \leq \rho \right\} \quad (220)$$

For easy reference, we give such generally nonconvex sets a name: *rank  $\rho$  subset of a positive semidefinite cone*. For  $\rho < M$  this subset, nonconvex for  $M > 1$ , resides on the positive semidefinite cone boundary.

##### 2.9.2.1.1 Exercise. Closure and rank $\rho$ subset.

Prove equality in (220). ▼

For example,

$$\partial \mathbb{S}_+^M = \overline{\left\{ A \in \mathbb{S}_+^M \mid \text{rank } A = M-1 \right\}} = \left\{ A \in \mathbb{S}_+^M \mid \text{rank } A \leq M-1 \right\} \quad (221)$$

In  $\mathbb{S}^2$ , each and every ray on the boundary of the positive semidefinite cone in isomorphic  $\mathbb{R}^3$  corresponds to a symmetric rank-1 matrix (Figure 46), but that does not hold in any higher dimension.

#### 2.9.2.2 Subspace tangent to open rank $\rho$ subset

When the positive semidefinite cone subset in (220) is left unclosed as in

$$\mathcal{M}(\rho) \triangleq \left\{ A \in \mathbb{S}_+^N \mid \text{rank } A = \rho \right\} \quad (222)$$

then we can specify a subspace tangent to the positive semidefinite cone at a particular member of manifold  $\mathcal{M}(\rho)$ . Specifically, the subspace  $\mathcal{R}_{\mathcal{M}}$  tangent to manifold  $\mathcal{M}(\rho)$  at  $B \in \mathcal{M}(\rho)$  [212, §5, prop.1.1]

$$\mathcal{R}_{\mathcal{M}}(B) \triangleq \{XB + BX^T \mid X \in \mathbb{R}^{N \times N}\} \subseteq \mathbb{S}^N \quad (223)$$

has dimension

$$\dim \text{svec } \mathcal{R}_{\mathcal{M}}(B) = \rho \left( N - \frac{\rho-1}{2} \right) = \rho(N-\rho) + \frac{\rho(\rho+1)}{2} \quad (224)$$

Tangent subspace  $\mathcal{R}_{\mathcal{M}}$  contains no member of the positive semidefinite cone  $\mathbb{S}_+^N$  whose rank exceeds  $\rho$ .

Subspace  $\mathcal{R}_{\mathcal{M}}(B)$  is a hyperplane supporting  $\mathbb{S}_+^N$  when  $B \in \mathcal{M}(N-1)$ . Another good example of tangent subspace is given in §E.7.2.0.2 by (2198);  $\mathcal{R}_{\mathcal{M}}(\mathbf{1}\mathbf{1}^T) = \mathbb{S}_c^{N\perp}$ , orthogonal complement to the *geometric center subspace*. (Figure 173 p.444)

### 2.9.2.3 Faces of PSD cone, their dimension *versus* rank

Define  $\mathcal{F}(\mathbb{S}_+^M \ni A)$  (174) as the smallest face, that contains a given positive semidefinite matrix  $A$ , of positive semidefinite cone  $\mathbb{S}_+^M$ . Then matrix  $A$ , having rank  $\rho$  and ordered diagonalization  $A = Q\Lambda Q^T \in \mathbb{S}_+^M$  (§A.5.1), is relatively interior to<sup>2.42</sup> [269] [27, §II.12] [126, §31.5.3] [268, §2.4]

$$\begin{aligned} \mathcal{F}(\mathbb{S}_+^M \ni A) &= \{X \in \mathbb{S}_+^M \mid \mathcal{N}(X) \supseteq \mathcal{N}(A)\} \\ &= \{X \in \mathbb{S}_+^M \mid \langle Q(I - \Lambda\Lambda^\dagger)Q^T, X \rangle = 0\} \\ &= \{Q\Lambda\Lambda^\dagger\Psi\Lambda\Lambda^\dagger Q^T \mid \Psi \in \mathbb{S}_+^M\} \\ &= Q\Lambda\Lambda^\dagger \mathbb{S}_+^M \Lambda\Lambda^\dagger Q^T \\ &= Q(:, 1:\rho) \mathbb{S}_+^\rho Q(:, 1:\rho)^T \\ &\simeq \mathbb{S}_+^\rho \end{aligned} \quad (225)$$

which is isomorphic with convex cone  $\mathbb{S}_+^\rho$ ; e.g.,  $Q \mathbb{S}_+^M Q^T = \mathbb{S}_+^M$ . The larger the nullspace of  $A$ , the smaller the face. (142) Thus dimension of the smallest face that contains given matrix  $A$  is

$$\dim \mathcal{F}(\mathbb{S}_+^M \ni A) = \rho(\rho + 1)/2 \quad (226)$$

in isomorphic  $\mathbb{R}^{M(M+1)/2}$ .

Each and every face of  $\mathbb{S}_+^M$  is isomorphic with a positive semidefinite cone having dimension the same as the face. Observe: not all dimensions are represented, and the only zero-dimensional face is the origin; e.g., a positive semidefinite cone has no facets:

#### 2.9.2.3.1 Table: Rank $\rho$ *versus* dimension of $\mathbb{S}_+^3$ faces

	$\rho$	$\dim \mathcal{F}(\mathbb{S}_+^3 \ni \text{rank-}\rho \text{ matrix})$
boundary	0	0
	$\leq 1$	1
	$\leq 2$	3
interior	$\leq 3$	6

For positive semidefinite cone  $\mathbb{S}_+^2$  in isometrically isomorphic  $\mathbb{R}^3$  depicted in Figure 46, rank-2 matrices belong to the interior of that face having dimension 3 (the entire closed cone), rank-1 matrices belong to relative interior of a face having dimension<sup>2.43</sup> 1, and the only rank-0 matrix is the point at the origin (the zero-dimensional face).

<sup>2.42</sup>For  $X \in \mathbb{S}_+^M$ ,  $A = Q\Lambda Q^T \in \mathbb{S}_+^M$ , show:  $\mathcal{N}(X) \supseteq \mathcal{N}(A) \Leftrightarrow \langle Q(I - \Lambda\Lambda^\dagger)Q^T, X \rangle = 0$ .

Given  $\langle Q(I - \Lambda\Lambda^\dagger)Q^T, X \rangle = 0 \Leftrightarrow \mathcal{R}(X) \perp \mathcal{N}(A)$ . (§A.7.4)

( $\Rightarrow$ ) Assume  $\mathcal{N}(X) \supseteq \mathcal{N}(A)$ , then  $\mathcal{R}(X) \perp \mathcal{N}(X) \supseteq \mathcal{N}(A)$ .

( $\Leftarrow$ ) Assume  $\mathcal{R}(X) \perp \mathcal{N}(A)$ , then  $X Q(I - \Lambda\Lambda^\dagger)Q^T = \mathbf{0} \Rightarrow \mathcal{N}(X) \supseteq \mathcal{N}(A)$ .  $\blacklozenge$

<sup>2.43</sup>The boundary constitutes all the one-dimensional faces, in  $\mathbb{R}^3$ , which are rays emanating from the origin.

### 2.9.2.3.2 Exercise. Bijective isometry.

Prove that the smallest face of positive semidefinite cone  $\mathbb{S}_+^M$ , containing a particular full-rank matrix  $A$  having ordered diagonalization  $Q\Lambda Q^T$ , is the entire cone: *id est*, prove  $Q\mathbb{S}_+^M Q^T = \mathbb{S}_+^M$  from (225).  $\blacktriangledown$

### 2.9.2.4 rank- $\rho$ face of PSD cone

Because each and every face of the positive semidefinite cone contains the origin (§2.8.0.0.1), each face belongs to a subspace of dimension the same as the face; from (225)

$$\mathcal{F}(\mathbb{S}_+^M \ni A) \subseteq \mathcal{S}_{\mathcal{F}} \triangleq Q(:, 1:\rho)\mathbb{S}^\rho Q(:, 1:\rho)^T \simeq \mathbb{S}^\rho \quad (227)$$

Because  $Q(:, 1:\rho)^T \mathbb{S}^\rho Q(:, 1:\rho) = \mathbb{S}^\rho$ , a surjection, projection of any matrix  $Y \in \mathbb{S}^M$  on this subspace  $\mathcal{S}_{\mathcal{F}} \subseteq \mathbb{S}^M$  is expressed  $P Y P$  (§E.7) where  $P \triangleq Q(:, 1:\rho)Q(:, 1:\rho)^T$ .  $\mathcal{S}_{\mathcal{F}}$  is the smallest subspace containing  $\mathcal{F}$ .

Each and every face of the positive semidefinite cone, having dimension less than that of the cone, is exposed. [276, §6] [242, §2.3.4] Any rank- $\rho < M$  positive semidefinite matrix  $A$  belongs to a face, of positive semidefinite cone  $\mathbb{S}_+^M$ , described by intersection with a hyperplane: for ordered diagonalization of  $A = Q\Lambda Q^T \in \mathbb{S}_+^M \ni \text{rank}(A) = \rho < M$

$$\begin{aligned} \mathcal{F}(\mathbb{S}_+^M \ni A) &= \{X \in \mathbb{S}_+^M \mid \langle Q(I - \Lambda\Lambda^\dagger)Q^T, X \rangle = 0\} \\ &= \left\{X \in \mathbb{S}_+^M \mid \left\langle Q\left(I - \begin{bmatrix} I \in \mathbb{S}^\rho & \mathbf{0} \\ \mathbf{0}^T & \mathbf{0} \end{bmatrix}\right)Q^T, X \right\rangle = 0\right\} \\ &= \mathbb{S}_+^M \cap \partial\mathcal{H}_+ \\ &\simeq \mathbb{S}_+^\rho \end{aligned} \quad (228)$$

Faces are doubly indexed: continuously indexed by orthogonal matrix  $Q$ , and discretely indexed by rank  $\rho$ . Each and every orthogonal matrix  $Q$  makes projectors  $Q(:, \rho+1:M)Q(:, \rho+1:M)^T$  indexed by  $\rho$ , in other words, each projector describing a normal<sup>2.44</sup> svec( $Q(:, \rho+1:M)Q(:, \rho+1:M)^T$ ) to a supporting hyperplane  $\partial\mathcal{H}_+$  (containing the origin) exposing a face (§2.11) of the positive semidefinite cone that contains rank- $\rho$  (and less) matrices.

### 2.9.2.4.1 Exercise. Simultaneously diagonalizable means commutative.

Given diagonalization of rank- $\rho \leq M$  positive semidefinite matrix  $A = Q\Lambda Q^T$  and any particular  $\Psi \succeq 0$ , both in  $\mathbb{S}^M$  from (225), show how  $I - \Lambda\Lambda^\dagger$  and  $\Lambda\Lambda^\dagger\Psi\Lambda\Lambda^\dagger$  share a complete set of eigenvectors.  $\blacktriangledown$

### 2.9.2.5 PSD cone face containing principal submatrix

A *principal submatrix* of a matrix  $A \in \mathbb{R}^{M \times M}$  is formed by discarding any particular subset of its rows and columns having the same indices. There are  $M!/(1!(M-1)!)$  principal  $1 \times 1$  submatrices,  $M!/(2!(M-2)!)$  principal  $2 \times 2$  submatrices, and so on, totaling  $2^M - 1$  principal submatrices including  $A$  itself. Principal submatrices of a symmetric matrix are symmetric. A given symmetric matrix has rank  $\rho$  iff it has a nonsingular  $\rho \times \rho$  principal submatrix but none larger. [332, §5-10] By loading vector  $y$  in test  $y^T A y$  (§A.2) with various binary patterns, it follows that any principal submatrix must be positive (semi)definite whenever  $A$  is (Theorem A.3.1.0.4). If positive semidefinite matrix  $A \in \mathbb{S}_+^M$  has principal submatrix of dimension  $\rho$  with rank  $r$ , then  $\text{rank } A \leq M - \rho + r$  by (1688).

<sup>2.44</sup>Any vectorized nonzero matrix  $\in \mathbb{S}_+^M$  is normal to a hyperplane supporting  $\mathbb{S}_+^M$  (§2.13.1) because PSD cone is selfdual. Normal on boundary exposes nonzero face by (333) (334).

Because each and every principal submatrix of a positive semidefinite matrix in  $\mathbb{S}_+^M$  is PSD, then each principal submatrix belongs to a certain face of cone  $\mathbb{S}_+^M$  by (226). Of special interest are full-rank positive semidefinite principal submatrices, for then description of smallest face becomes simpler. We can find the smallest face, that contains a particular complete full-rank principal submatrix of  $A$ , by embedding that submatrix in a  $\mathbf{0}$  matrix of the same dimension as  $A$ : Were  $\Phi$  a binary diagonal matrix

$$\Phi = \delta^2(\Phi) \in \mathbb{S}^M, \quad \Phi_{ii} \in \{0, 1\} \quad (229)$$

having diagonal entry 0 corresponding to a discarded row and column from  $A \in \mathbb{S}_+^M$ , then any principal submatrix <sup>2.45</sup> so embedded can be expressed  $\Phi A \Phi$ ; *id est*, for an embedded principal submatrix  $\Phi A \Phi \in \mathbb{S}_+^M \ni \text{rank } \Phi A \Phi = \text{rank } \Phi \leq \text{rank } A$

$$\begin{aligned} \mathcal{F}\left(\mathbb{S}_+^M \ni \Phi A \Phi\right) &= \{X \in \mathbb{S}_+^M \mid \mathcal{N}(X) \supseteq \mathcal{N}(\Phi A \Phi)\} \\ &= \{X \in \mathbb{S}_+^M \mid \langle I - \Phi, X \rangle = 0\} \\ &= \{\Phi \Psi \Phi \mid \Psi \in \mathbb{S}_+^M\} \\ &\simeq \mathbb{S}_+^{\text{rank } \Phi} \end{aligned} \quad (230)$$

The smallest face that contains an embedded principal submatrix, whose rank is not necessarily full, may be expressed like (225): For embedded principal submatrix  $\Phi A \Phi \in \mathbb{S}_+^M \ni \text{rank } \Phi A \Phi \leq \text{rank } \Phi$ , apply ordered diagonalization instead to

$$\hat{\Phi}^T A \hat{\Phi} = U \Upsilon U^T \in \mathbb{S}_+^{\text{rank } \Phi} \quad (231)$$

where  $U^{-1} = U^T$  is an orthogonal matrix and  $\Upsilon = \delta^2(\Upsilon)$  is diagonal. Then

$$\begin{aligned} \mathcal{F}\left(\mathbb{S}_+^M \ni \Phi A \Phi\right) &= \{X \in \mathbb{S}_+^M \mid \mathcal{N}(X) \supseteq \mathcal{N}(\Phi A \Phi)\} \\ &= \{X \in \mathbb{S}_+^M \mid \langle \hat{\Phi} U(I - \Upsilon \Upsilon^\dagger) U^T \hat{\Phi}^T + I - \Phi, X \rangle = 0\} \\ &= \{\hat{\Phi} U \Upsilon \Upsilon^\dagger \Psi \Upsilon \Upsilon^\dagger U^T \hat{\Phi}^T \mid \Psi \in \mathbb{S}_+^{\text{rank } \Phi}\} \\ &\simeq \mathbb{S}_+^{\text{rank } \Phi A \Phi} \end{aligned} \quad (232)$$

where binary diagonal matrix  $\Phi$  is partitioned into nonzero and zero columns by permutation  $\Xi \in \mathbb{R}^{M \times M}$ ;

$$\Phi \Xi^T \triangleq [\hat{\Phi} \ \mathbf{0}] \in \mathbb{R}^{M \times M}, \quad \text{rank } \hat{\Phi} = \text{rank } \Phi, \quad \Phi = \hat{\Phi} \hat{\Phi}^T \in \mathbb{S}^M, \quad \hat{\Phi}^T \hat{\Phi} = I \quad (233)$$

Any embedded principal submatrix may be expressed

$$\Phi A \Phi = \hat{\Phi} \hat{\Phi}^T A \hat{\Phi} \hat{\Phi}^T \in \mathbb{S}_+^M \quad (234)$$

where  $\hat{\Phi}^T A \hat{\Phi} \in \mathbb{S}_+^{\text{rank } \Phi}$  extracts the principal submatrix whereas  $\hat{\Phi} \hat{\Phi}^T A \hat{\Phi} \hat{\Phi}^T$  embeds it.

#### 2.9.2.5.1 Example. Smallest face containing disparate elements.

Smallest face formula (225) can be altered to accommodate a union of points  $\{A_i \in \mathbb{S}_+^M\}$ :

$$\mathcal{F}\left(\mathbb{S}_+^M \supset \bigcup_i A_i\right) = \left\{X \in \mathbb{S}_+^M \mid \mathcal{N}(X) \supseteq \bigcap_i \mathcal{N}(A_i)\right\} \quad (235)$$

---

<sup>2.45</sup>To express a leading principal submatrix, for example,  $\Phi = \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0}^T & \mathbf{0} \end{bmatrix}$ .

To see that, imagine two vectorized matrices  $A_1$  and  $A_2$  on diametrically opposed sides of the positive semidefinite cone  $\mathbb{S}_+^2$  boundary pictured in Figure 46. Regard svec  $A_1$  as normal to a hyperplane in  $\mathbb{R}^3$  containing a vectorized basis for its nullspace: svec basis  $\mathcal{N}(A_1)$  (§2.5.2.0.1). Similarly, there is a second hyperplane containing svec basis  $\mathcal{N}(A_2)$  having normal svec  $A_2$ . While each hyperplane is two-dimensional, each nullspace has only one affine dimension because  $A_1$  and  $A_2$  are rank-1. Because our interest is only that part of the nullspace in the positive semidefinite cone, then by

$$\langle X, A_i \rangle = 0 \Leftrightarrow X A_i = A_i X = \mathbf{0}, \quad X, A_i \in \mathbb{S}_+^M \quad (1748)$$

we may ignore the fact that vectorized nullspace svec basis  $\mathcal{N}(A_i)$  is a proper subspace of the hyperplane. We may think, instead, in terms of whole hyperplanes because equivalence (1748) says that the positive semidefinite cone effectively filters that subset of the hyperplane, whose normal is  $A_i$ , constituting  $\mathcal{N}(A_i) = \mathcal{N}(A_i^T)$ .

And so hyperplane intersection makes a line intersecting positive semidefinite cone  $\mathbb{S}_+^2$  but only at the origin. In this hypothetical example, the smallest face containing those two matrices  $A_1$  and  $A_2$  therefore comprises the entire cone because every positive semidefinite matrix has nullspace containing  $\mathbf{0}$ . The smaller the intersection, the larger the smallest face.  $\square$

#### 2.9.2.5.2 Exercise. Disparate elements.

Prove that (235) holds for an arbitrary set  $\{A_i \in \mathbb{S}_+^M \forall i \in \mathcal{I}\}$ . One way is by showing  $\bigcap \mathcal{N}(A_i) \cap \mathbb{S}_+^M = \text{conv}(\{A_i\})^\perp \cap \mathbb{S}_+^M$ ; with perpendicularity  $^\perp$  as in (377).<sup>2.46</sup>  $\blacktriangledown$

#### 2.9.2.6 face of all PSD matrices having same principal submatrix

Now we ask what is the smallest face of the positive semidefinite cone containing all matrices having a complete principal submatrix in common; in other words, that face containing all PSD matrices (of any rank) with particular entries fixed - the smallest face containing all PSD matrices whose fixed entries correspond to some given embedded principal submatrix  $\Phi A \Phi$ . To maintain generality,<sup>2.47</sup> we move an extracted principal submatrix  $\hat{\Phi}^T A \hat{\Phi} \in \mathbb{S}_+^{\text{rank } \Phi}$  into leading position via permutation  $\Xi$  from (233): for  $A \in \mathbb{S}_+^M$

$$\Xi A \Xi^T \triangleq \begin{bmatrix} \hat{\Phi}^T A \hat{\Phi} & B \\ B^T & C \end{bmatrix} \in \mathbb{S}_+^M \quad (236)$$

By properties of partitioned PSD matrices in §A.4.0.1,

$$\text{basis } \mathcal{N} \left( \begin{bmatrix} \hat{\Phi}^T A \hat{\Phi} & B \\ B^T & C \end{bmatrix} \right) \supseteq \begin{bmatrix} \mathbf{0} \\ I - CC^\dagger \end{bmatrix} \quad (237)$$

Hence  $\mathcal{N}(\Xi A \Xi^T) \supseteq \mathcal{N}(\Xi A \Xi^T) \not\supseteq \text{span} \begin{bmatrix} \mathbf{0} \\ I \end{bmatrix}$  in a smallest face  $\mathcal{F}$  formula<sup>2.48</sup> because all PSD matrices, given fixed principal submatrix, are admitted: Define a set of all PSD matrices

$$\mathcal{S}_+ \triangleq \left\{ A = \Xi^T \begin{bmatrix} \hat{\Phi}^T A \hat{\Phi} & B \\ B^T & C \end{bmatrix} \Xi \succeq 0 \mid B \in \mathbb{R}^{\text{rank } \Phi \times M - \text{rank } \Phi}, C \in \mathbb{S}_+^{M - \text{rank } \Phi} \right\} \quad (238)$$

<sup>2.46</sup> Hint: (1748) (2095).

<sup>2.47</sup> to fix any principal submatrix; not only leading principal submatrices.

<sup>2.48</sup> meaning, more pertinently,  $I - \Phi$  is dropped from (232).

having fixed embedded principal submatrix  $\Phi A \Phi = \Xi^T \begin{bmatrix} \hat{\Phi}^T A \hat{\Phi} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{0} \end{bmatrix} \Xi$ . So

$$\begin{aligned} \mathcal{F}\left(\mathbb{S}_+^M \supseteq \mathcal{S}_+\right) &= \left\{ X \in \mathbb{S}_+^M \mid \mathcal{N}(X) \supseteq \mathcal{N}(\mathcal{S}_+) \right\} \\ &= \left\{ X \in \mathbb{S}_+^M \mid \langle \hat{\Phi} U(I - \Upsilon \Upsilon^\dagger) U^T \hat{\Phi}^T, X \rangle = 0 \right\} \\ &= \left\{ \Xi^T \begin{bmatrix} U \Upsilon \Upsilon^\dagger & \mathbf{0} \\ \mathbf{0}^T & I \end{bmatrix} \Psi \begin{bmatrix} \Upsilon \Upsilon^\dagger U^T & \mathbf{0} \\ \mathbf{0}^T & I \end{bmatrix} \Xi \mid \Psi \in \mathbb{S}_+^M \right\} \\ &\simeq \mathbb{S}_+^{M - \text{rank } \Phi + \text{rank } \Phi A \Phi} \end{aligned} \tag{239}$$

$\Xi = I$  whenever  $\Phi A \Phi$  denotes a leading principal submatrix. Smallest face of the positive semidefinite cone, containing all matrices having the same full-rank principal submatrix ( $\Upsilon \Upsilon^\dagger = I$ ,  $\Upsilon \succeq 0$ ), is the entire cone (Exercise 2.9.2.3.2).

#### 2.9.2.7 Extreme directions of positive semidefinite cone

Because the positive semidefinite cone is pointed (§2.7.2.1.2), there is a one-to-one correspondence of one-dimensional faces with extreme directions in any dimension  $M$ ; *id est*, because of the *cone faces lemma* (§2.8.0.0.1) and direct correspondence of exposed faces to faces of  $\mathbb{S}_+^M$ , it follows: there is no one-dimensional face of the positive semidefinite cone that is not a ray emanating from the origin.

Symmetric dyads constitute the set of all extreme directions: For  $M > 1$

$$\{yy^T \in \mathbb{S}^M \mid y \in \mathbb{R}^M\} \subset \partial \mathbb{S}_+^M \tag{240}$$

this superset of extreme directions (infinite in number, *confer*(190)) for the positive semidefinite cone is a proper subset of the boundary when  $M > 2$ . By *extremes theorem* 2.8.1.1.1, the convex hull of extreme rays and origin is the positive semidefinite cone: (§2.8.1.2.1)

$$\text{conv}\{yy^T \in \mathbb{S}^M \mid y \in \mathbb{R}^M\} = \left\{ \sum_{i=1}^{\infty} b_i z_i z_i^T \mid z_i \in \mathbb{R}^M, b \succeq 0 \right\} = \mathbb{S}_+^M \tag{241}$$

For two-dimensional matrices ( $M = 2$ , Figure 46)

$$\{yy^T \in \mathbb{S}^2 \mid y \in \mathbb{R}^2\} = \partial \mathbb{S}_+^2 \tag{242}$$

while for one-dimensional matrices, in exception, ( $M = 1$ , §2.7)

$$\{yy \in \mathbb{S} \mid y \neq \mathbf{0}\} = \text{intr } \mathbb{S}_+ \tag{243}$$

Each and every extreme direction  $yy^T$  makes the same angle with the Identity matrix in isomorphic  $\mathbb{R}^{M(M+1)/2}$ , dependent only on dimension; *videlicet*, 2.49

$$\angle(yyyy^T, I) = \arccos \frac{\langle yyyy^T, I \rangle}{\|yyyy^T\|_F \|I\|_F} = \arccos \left( \frac{1}{\sqrt{M}} \right) \quad \forall y \in \mathbb{R}^M \tag{244}$$

This means the positive semidefinite cone broadens in higher dimension.

---

**2.49** Analogy with respect to the *EDM cone* is considered in [211, p.162] where it is found: angle is not constant. Extreme directions of the EDM cone can be found in §6.4.3.2. The cone's axis is  $-E = \mathbf{1}\mathbf{1}^T - I$  (1240).

### 2.9.2.7.1 Example. Positive semidefinite matrix from extreme directions.

Diagonalizability (§A.5) of symmetric matrices yields the following results:

Any positive semidefinite matrix (1604) in  $\mathbb{S}^M$  can be written in the form

$$A = \sum_{i=1}^M \lambda_i z_i z_i^T = \hat{A} \hat{A}^T = \sum_i \hat{a}_i \hat{a}_i^T \succeq 0, \quad \lambda \succeq 0 \quad (245)$$

a conic combination of linearly independent extreme directions ( $\hat{a}_i \hat{a}_i^T$  or  $z_i z_i^T$ ,  $\|z_i\|=1$ ) of the positive semidefinite cone, where  $\lambda$  is a vector of eigenvalues.

If we limit consideration to all symmetric positive semidefinite matrices bounded via unity trace

$$\mathcal{C} \triangleq \{A \succeq 0 \mid \text{tr } A = 1\} \quad (92)$$

then any matrix  $A$  from that set may be expressed as a convex combination of linearly independent extreme directions;

$$A = \sum_{i=1}^M \lambda_i z_i z_i^T \in \mathcal{C}, \quad \mathbf{1}^T \lambda = 1, \quad \lambda \succeq 0 \quad (246)$$

Implications are:

1. set  $\mathcal{C}$  is convex (an intersection of PSD cone with hyperplane),
2. because the set of eigenvalues corresponding to a given square matrix  $A$  is unique (§A.5.0.1), no single eigenvalue can exceed 1; *id est*,  $I \succeq A$
3. and the converse holds: set  $\mathcal{C}$  is an instance of Fantope (92). □

### 2.9.2.7.2 Exercise. Extreme directions of positive semidefinite cone.

Prove, directly from definition (189), that symmetric dyads (240) constitute the set of all extreme directions of the positive semidefinite cone. ▼

### 2.9.2.8 Positive semidefinite cone is generally not circular

Extreme angle equation (244) suggests that the positive semidefinite cone might be invariant to rotation about its axis of revolution; *id est*, a circular cone. We investigate this now:

#### 2.9.2.8.1 Definition. Circular cone:<sup>2.50</sup>

a pointed closed convex cone having hyperspherical sections orthogonal to its *axis of revolution* about which the cone is invariant to rotation. △

A *conic section* is the intersection of a cone with any hyperplane. In three dimensions, an intersecting plane perpendicular to a circular cone's axis of revolution produces a section bounded by a circle. (Figure 49) A prominent example of circular cone, in convex analysis, is Lorentz cone (181). We also find that the positive semidefinite cone and cone of Euclidean distance matrices are circular cones, but only in low dimension.

The positive semidefinite cone has axis of revolution that is the ray (base  $\mathbf{0}$ ) through the Identity matrix  $I$ . Consider a set of normalized extreme directions of the positive semidefinite cone: for some arbitrary positive constant  $a \in \mathbb{R}_+$

$$\{yy^T \in \mathbb{S}^M \mid \|y\| = \sqrt{a}\} \subset \partial \mathbb{S}_+^M \quad (247)$$

---

<sup>2.50</sup>A circular cone is assumed convex throughout, although not so by other authors. We also assume a *right* circular cone.

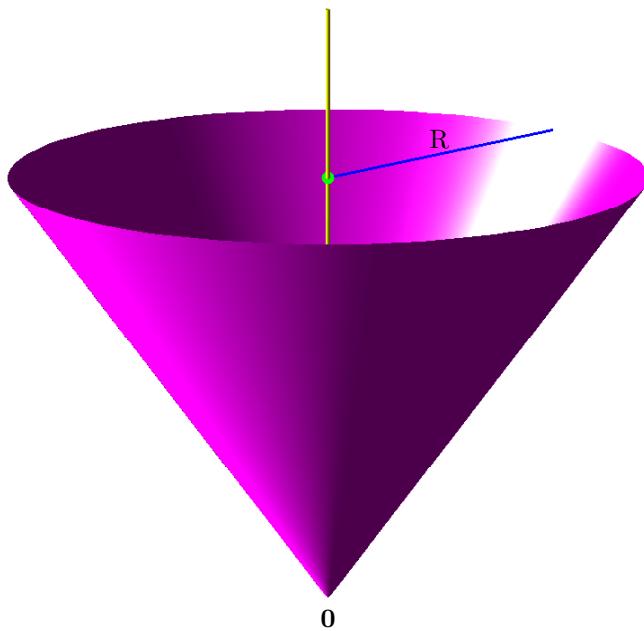


Figure 49: This solid circular cone in  $\mathbb{R}^3$  continues upward infinitely. Axis of revolution is illustrated as vertical line through origin.  $R$  represents radius: distance measured from an extreme direction to axis of revolution. Were this a Lorentz cone, any plane slice containing axis of revolution would make a right angle.

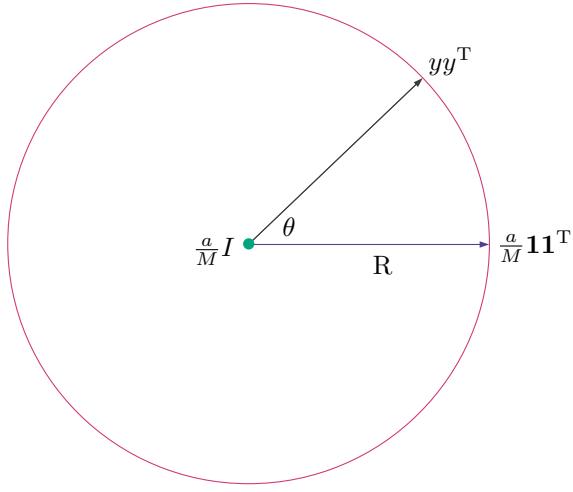


Figure 50: Illustrated is a section, perpendicular to axis of revolution, of circular cone from Figure 49. Radius  $R$  is distance from any extreme direction to axis at  $\frac{a}{M}I$ . Vector  $\frac{a}{M}\mathbf{1}\mathbf{1}^T$  is an arbitrary reference by which to measure angle  $\theta$ .

The distance from each extreme direction to the axis of revolution is radius

$$R \triangleq \inf_c \|yy^T - cI\|_F = a\sqrt{1 - \frac{1}{M}} \quad (248)$$

which is the distance from  $yy^T$  to  $\frac{a}{M}I$ ; the length of vector  $yy^T - \frac{a}{M}I$ .

Because distance  $R$  (in a particular dimension) from the axis of revolution to each and every normalized extreme direction is identical, the extreme directions lie on the boundary of a hypersphere in isometrically isomorphic  $\mathbb{R}^{M(M+1)/2}$ . From Example 2.9.2.7.1, the convex hull (excluding vertex at the origin) of the normalized extreme directions is a conic section

$$\mathcal{C} \triangleq \text{conv}\{yy^T \mid y \in \mathbb{R}^M, y^T y = a\} = \mathbb{S}_+^M \cap \{A \in \mathbb{S}^M \mid \langle I, A \rangle = a\} \quad (249)$$

orthogonal to Identity matrix  $I$ ;

$$\left\langle \mathcal{C} - \frac{a}{M}I, I \right\rangle = \text{tr}(\mathcal{C} - \frac{a}{M}I) = 0 \quad (250)$$

**Proof.** Although the positive semidefinite cone possesses some characteristics of a circular cone, we can show it is not by demonstrating shortage of extreme directions; *id est*, some extreme directions corresponding to each and every angle of rotation about the axis of revolution are nonexistent: Referring to Figure 50, [442, §1-7]

$$\cos \theta = \frac{\langle \frac{a}{M}\mathbf{1}\mathbf{1}^T - \frac{a}{M}I, yy^T - \frac{a}{M}I \rangle}{a^2(1 - \frac{1}{M})} \quad (251)$$

Solving for vector  $y$  we get

$$a(1 + (M-1)\cos \theta) = (\mathbf{1}^T y)^2 \quad (252)$$

which does not have real solution  $\forall 0 \leq \theta \leq 2\pi$  in every matrix dimension  $M$ .  $\diamond$

From the foregoing proof we can conclude that the positive semidefinite cone might be circular but only in matrix dimensions 1 and 2. Because of a shortage of extreme directions, conic section (249) cannot be hyperspherical by the *extremes theorem* (§2.8.1.1.1, Figure 45).

#### 2.9.2.8.2 Exercise. Circular semidefinite cone.

Prove the positive semidefinite cone to be circular in matrix dimensions 1 and 2 while it is a rotation of Lorentz cone (181) in matrix dimension 2. [2.51](#) ▼

#### 2.9.2.8.3 Example. Positive semidefinite cone inscription in three dimensions.

**Theorem.** *Geršgorin discs.*

[228, §6.1] [405] [285, p.140]

Given  $A = [A_{ij}] \in \mathbb{S}^m$ , all its eigenvalues belong to a union of  $m$  closed intervals on the real line; for  $p \in \mathbb{R}_+^m$

$$\lambda(A) \in \bigcup_{i=1}^m \left\{ \xi \in \mathbb{R} \mid |\xi - A_{ii}| \leq \varrho_i \triangleq \frac{1}{p_i} \sum_{\substack{j=1 \\ j \neq i}}^m p_j |A_{ij}| \right\} = \bigcup_{i=1}^m [A_{ii} - \varrho_i, A_{ii} + \varrho_i] \quad (253)$$

Furthermore, if a union of  $k$  of these  $m$  [intervals] forms a connected region that is disjoint from all the remaining  $n-k$  [intervals], then there are precisely  $k$  eigenvalues of  $A$  in this region. ◇

To apply the theorem to determine positive semidefiniteness of symmetric matrix  $A$ , observe that for each  $i$  we must have

$$A_{ii} \geq \varrho_i \quad (254)$$

Suppose

$$m = 2 \quad (255)$$

so  $A \in \mathbb{S}^2$ . Vectorizing  $A$  as in (57), svec  $A$  belongs to isometrically isomorphic  $\mathbb{R}^3$ . Then we have  $m2^{m-1}=4$  inequalities, in the matrix entries  $A_{ij}$  with Geršgorin parameters  $p = [p_i] \in \mathbb{R}_+^2$ ,

$$\begin{aligned} p_1 A_{11} &\geq \pm p_2 A_{12} \\ p_2 A_{22} &\geq \pm p_1 A_{12} \end{aligned} \quad (256)$$

which describe an intersection of four halfspaces in  $\mathbb{R}^{m(m+1)/2}$ . That intersection creates the proper polyhedral cone  $\mathcal{K}$  (§2.12.1) whose construction is illustrated in Figure 51. Drawn truncated is the boundary of the positive semidefinite cone  $\text{svec } \mathbb{S}_+^2$  and the bounding hyperplanes supporting  $\mathcal{K}$ .

Created by means of Geršgorin discs,  $\mathcal{K}$  always belongs to the positive semidefinite cone for any nonnegative value of  $p \in \mathbb{R}_+^m$ . Hence any point in  $\mathcal{K}$  corresponds to some positive semidefinite matrix  $A$ . Only the extreme directions of  $\mathcal{K}$  intersect the positive semidefinite cone boundary in this dimension; the four extreme directions of  $\mathcal{K}$  are extreme directions of the positive semidefinite cone. As  $p_1/p_2$  increases in value from 0, two extreme directions of  $\mathcal{K}$  sweep the entire boundary of this positive semidefinite cone. Because the entire positive semidefinite cone can be swept by  $\mathcal{K}$ , the system of linear inequalities

$$Y^T \text{svec } A \triangleq \begin{bmatrix} p_1 & \pm p_2/\sqrt{2} & 0 \\ 0 & \pm p_1/\sqrt{2} & p_2 \end{bmatrix} \text{svec } A \succeq 0 \quad (257)$$

---

[2.51](#) Hint: Given cone  $\left\{ \begin{bmatrix} \alpha & \beta/\sqrt{2} \\ \beta/\sqrt{2} & \gamma \end{bmatrix} \mid \sqrt{\alpha^2 + \beta^2} \leq \gamma \right\}$ , show  $\frac{1}{\sqrt{2}} \begin{bmatrix} \gamma + \alpha & \beta \\ \beta & \gamma - \alpha \end{bmatrix}$  to be a vector rotation that is positive semidefinite under the same inequality.

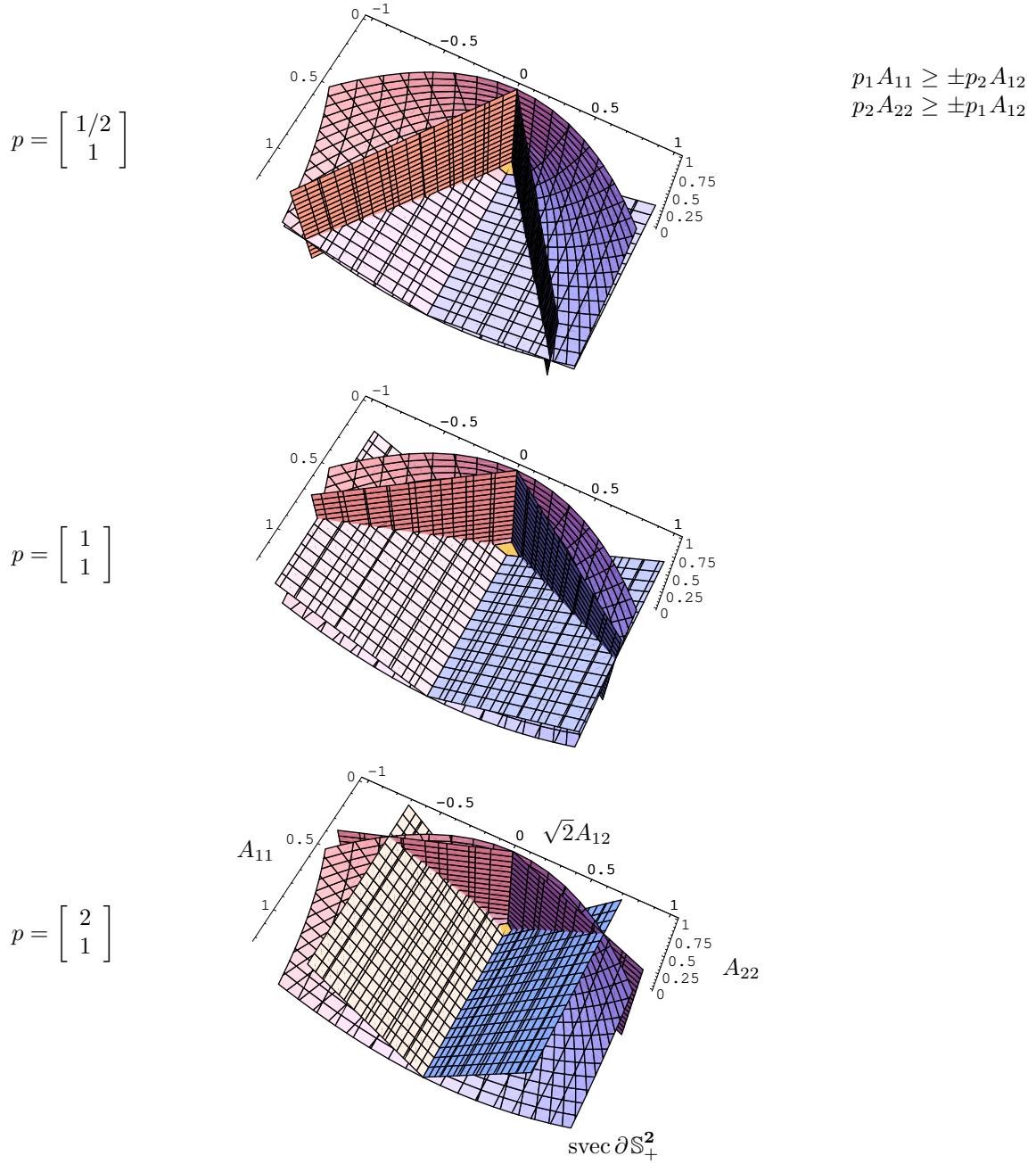


Figure 51: Proper polyhedral cone  $\mathcal{K}$ , created by intersection of halfspaces, inscribes PSD cone in isometrically isomorphic  $\mathbb{R}^3$  as predicted by *Geršgorin discs theorem* for  $A = [A_{ij}] \in \mathbb{S}^2$ . Hyperplanes supporting  $\mathcal{K}$  intersect along boundary of PSD cone. Four extreme directions of  $\mathcal{K}$  coincide with extreme directions of PSD cone.

(when made dynamic) can replace a semidefinite constraint  $A \succeq 0$ ; *id est*, for

$$\mathcal{K} = \{z \mid Y^T z \succeq 0\} \subset \text{svec } \mathbb{S}_+^m \quad (258)$$

given  $p$  where  $Y \in \mathbb{R}^{m(m+1)/2 \times m2^{m-1}}$

$$\text{svec } A \in \mathcal{K} \Rightarrow A \in \mathbb{S}_+^m \quad (259)$$

but

$$\exists p \ni Y^T \text{svec } A \succeq 0 \Leftrightarrow A \succeq 0 \quad (260)$$

In other words, *diagonal dominance* [228, p.349, §7.2.3]

$$A_{ii} \geq \sum_{\substack{j=1 \\ j \neq i}}^m |A_{ij}|, \quad \forall i = 1 \dots m \quad (261)$$

is generally only a sufficient condition for membership to the PSD cone. But by dynamic weighting  $p$  in this dimension, diagonal dominance was made necessary and sufficient.

□

In higher dimension ( $m > 2$ ), boundary of the positive semidefinite cone is no longer constituted completely by its extreme directions (symmetric rank-one matrices); its geometry becomes intricate. How all the extreme directions can be swept by an inscribed polyhedral cone,<sup>2.52</sup> similarly to the foregoing example, remains an open question.

#### 2.9.2.8.4 Exercise. Dual inscription.

Find dual proper polyhedral cone  $\mathcal{K}^*$  from Figure 51. ▼

#### 2.9.2.9 Boundary constituents of the positive semidefinite cone

**2.9.2.9.1 Lemma.** *Sum of positive semidefinite matrices.* (confer (1620))  
For  $A, B \in \mathbb{S}_+^M$

$$\text{rank}(A + B) = \text{rank}(\mu A + (1 - \mu)B) \quad (262)$$

over open interval  $(0, 1)$  of  $\mu$ . ◇

**Proof.** Any positive semidefinite matrix belonging to the PSD cone has an eigenvalue decomposition that is a positively scaled sum of linearly independent symmetric dyads. By the *linearly independent dyads definition* in §B.1.1.0.1, rank of the sum  $A+B$  is equivalent to the number of linearly independent dyads constituting it. Linear independence is insensitive to further positive scaling by  $\mu$ . The assumption of positive semidefiniteness prevents annihilation of any dyad from the sum  $A+B$ . ♦

**2.9.2.9.2 Example.** *Rank function quasiconcavity.* (confer §3.14)  
For  $A, B \in \mathbb{R}^{m \times n}$  [228, §0.4]

$$\text{rank } A + \text{rank } B \geq \text{rank}(A + B) \quad (263)$$

that follows from the fact [368, §3.6]

$$\dim \mathcal{R}(A) + \dim \mathcal{R}(B) = \dim \mathcal{R}(A + B) + \dim (\mathcal{R}(A) \cap \mathcal{R}(B)) \quad (264)$$

---

<sup>2.52</sup>It is not necessary to sweep the entire boundary in higher dimension.

For  $A, B \in \mathbb{S}_+^M$

$$\text{rank } A + \text{rank } B \geq \text{rank}(A + B) \geq \min\{\text{rank } A, \text{rank } B\} \quad (1620)$$

that follows from the fact

$$\mathcal{N}(A + B) = \mathcal{N}(A) \cap \mathcal{N}(B), \quad A, B \in \mathbb{S}_+^M \quad (163)$$

Rank is a *quasiconcave* function on  $\mathbb{S}_+^M$  because the right-hand inequality in (1620) has the concave form (654); *videlicet*, Lemma 2.9.2.9.1.  $\square$

From this example we see, unlike convex functions, *quasiconvex* functions are not necessarily continuous. (§3.14) We also glean:

**2.9.2.9.3 Theorem.** *Convex subsets of positive semidefinite cone.*  
Subsets of the positive semidefinite cone  $\mathbb{S}_+^M$ , for  $0 \leq \rho \leq M$

$$\mathbb{S}_+^M(\rho) \triangleq \{X \in \mathbb{S}_+^M \mid \text{rank } X \geq \rho\} \quad (265)$$

are pointed convex cones, but not closed unless  $\rho = 0$ ; *id est*,  $\mathbb{S}_+^M(0) = \mathbb{S}_+^M$ .  $\diamond$

**Proof.** Given  $\rho$ , a subset  $\mathbb{S}_+^M(\rho)$  is convex if and only if convex combination of any two members has rank at least  $\rho$ . That is confirmed by applying identity (262) from Lemma 2.9.2.9.1 to (1620); *id est*, for  $A, B \in \mathbb{S}_+^M(\rho)$  on closed interval  $[0, 1]$  of  $\mu$

$$\text{rank}(\mu A + (1 - \mu)B) \geq \min\{\text{rank } A, \text{rank } B\} \quad (266)$$

It can similarly be shown, almost identically to proof of the lemma: any conic combination of  $A, B$  in subset  $\mathbb{S}_+^M(\rho)$  remains a member; *id est*,  $\forall \zeta, \xi \geq 0$

$$\text{rank}(\zeta A + \xi B) \geq \min\{\text{rank}(\zeta A), \text{rank}(\xi B)\} \quad (267)$$

Therefore,  $\mathbb{S}_+^M(\rho)$  is a convex cone.  $\blacklozenge$

Another proof of convexity can be made by projection arguments:

#### 2.9.2.10 Projection on $\mathbb{S}_+^M(\rho)$

Because these cones  $\mathbb{S}_+^M(\rho)$  indexed by  $\rho$  (265) are convex, projection on them is straightforward. Given a symmetric matrix  $H$  having diagonalization  $H \triangleq Q\Lambda Q^T \in \mathbb{S}^M$  (§A.5.1) with eigenvalues  $\Lambda$  arranged in nonincreasing order, then its *Euclidean projection* (minimum-distance projection) on  $\mathbb{S}_+^M(\rho)$

$$P_{\mathbb{S}_+^M(\rho)} H = Q \Upsilon^* Q^T \quad (268)$$

corresponds to a map of its eigenvalues:

$$\Upsilon_{ii}^* = \begin{cases} \max\{\epsilon, \Lambda_{ii}\}, & i=1 \dots \rho \\ \max\{0, \Lambda_{ii}\}, & i=\rho+1 \dots M \end{cases} \quad (269)$$

where  $\epsilon$  is positive but arbitrarily close to 0.

### 2.9.2.10.1 Exercise. Projection on open convex cones.

Prove (269) using Theorem E.9.2.0.1. ▼

Because each  $H \in \mathbb{S}^M$  has unique projection on  $\mathbb{S}_+^M(\rho)$  (despite possibility of repeated eigenvalues in  $\Lambda$ ), we may conclude it is a convex set by the *Bunt-Motzkin theorem* (§E.9.0.0.1).

Compare (269) to the well-known result regarding Euclidean projection on a rank  $\rho$  subset of the positive semidefinite cone (§2.9.2.1)

$$\mathbb{S}_+^M \setminus \mathbb{S}_+^M(\rho+1) = \{X \in \mathbb{S}_+^M \mid \text{rank } X \leq \rho\} \quad (220)$$

$$P_{\mathbb{S}_+^M \setminus \mathbb{S}_+^M(\rho+1)} H = Q \Upsilon^* Q^T \quad (270)$$

As proved in §7.1.4, this projection of  $H$  corresponds to the eigenvalue map

$$\Upsilon_{ii}^* = \begin{cases} \max \{0, \Lambda_{ii}\}, & i=1 \dots \rho \\ 0, & i=\rho+1 \dots M \end{cases} \quad (1488)$$

Together these two results (269) and (1488) mean: A higher-rank solution to projection on the positive semidefinite cone lies arbitrarily close to any given lower-rank projection, but not *vice versa*. Were the number of nonnegative eigenvalues in  $\Lambda$  known *a priori* not to exceed  $\rho$ , then these two different projections would produce identical results in the limit  $\epsilon \rightarrow 0$ .

### 2.9.2.11 Uniting constituents

Interior of the PSD cone  $\text{intr } \mathbb{S}_+^M$  is convex by Theorem 2.9.2.9.3, for example, because all positive semidefinite matrices having rank  $M$  constitute the cone interior.

All positive semidefinite matrices of rank less than  $M$  constitute the cone boundary; an amalgam of positive semidefinite matrices of different rank. Thus each nonconvex subset of positive semidefinite matrices, for  $0 < \rho < M$

$$\{Y \in \mathbb{S}_+^M \mid \text{rank } Y = \rho\} \quad (271)$$

having rank  $\rho$  successively 1 lower than  $M$ , appends a nonconvex constituent to the cone boundary; but only in their union is the boundary complete: (*confer* §2.9.2)

$$\partial \mathbb{S}_+^M = \bigcup_{\rho=0}^{M-1} \{Y \in \mathbb{S}_+^M \mid \text{rank } Y = \rho\} \quad (272)$$

The composite sequence, the cone interior in union with each successive constituent, remains convex at each step; *id est*, for  $0 \leq k \leq M$

$$\bigcup_{\rho=k}^M \{Y \in \mathbb{S}_+^M \mid \text{rank } Y = \rho\} \quad (273)$$

is convex for each  $k$  by Theorem 2.9.2.9.3.

### 2.9.2.12 Peeling constituents

Proceeding the other way: To peel constituents off the complete positive semidefinite cone boundary, one starts by removing the origin; the only rank-0 positive semidefinite matrix. What remains is convex. Next, the extreme directions are removed because they constitute all the rank-1 positive semidefinite matrices. What remains is again convex, and so on. Proceeding in this manner eventually removes the entire boundary leaving, at last, the convex interior of the PSD cone; all the positive definite matrices.

### 2.9.2.12.1 Exercise. Difference $A - B$ .

What about a difference of matrices  $A, B$  belonging to the PSD cone? Show:

- Difference of any two points on the boundary belongs to the boundary or exterior.
- Difference  $A - B$ , where  $A$  belongs to the boundary while  $B$  is interior, belongs to the exterior.  $\blacktriangledown$

### 2.9.3 Barvinok's proposition

Barvinok posits existence and quantifies an upper bound on rank of a positive semidefinite matrix belonging to the intersection of the PSD cone with an affine subset:

#### 2.9.3.0.1 Proposition. Affine intersection with PSD cone. [27, §II.13] [25, §2.2]

Consider finding a matrix  $X \in \mathbb{S}^N$  satisfying

$$X \succeq 0, \quad \langle A_j, X \rangle = b_j, \quad j=1 \dots m \quad (2272)$$

given nonzero linearly independent (vectorized)  $A_j \in \mathbb{S}^N$  and real  $b_j$ . Define the affine subset

$$\mathcal{A} \triangleq \{X \mid \langle A_j, X \rangle = b_j, j=1 \dots m\} \subseteq \mathbb{S}^N \quad (2273)$$

If the intersection  $\mathcal{A} \cap \mathbb{S}_+^N$  is nonempty given a number  $m$  of equalities, then there exists a matrix  $X \in \mathcal{A} \cap \mathbb{S}_+^N$  such that

$$\text{rank } X (\text{rank } X + 1)/2 \leq m \quad (274)$$

whence the upper bound<sup>2.53</sup>

$$\text{rank } X \leq \left\lfloor \frac{\sqrt{8m+1}-1}{2} \right\rfloor \quad (275)$$

Given desired rank instead, equivalently,

$$m < (\text{rank } X + 1)(\text{rank } X + 2)/2 \quad (276)$$

An extreme point of  $\mathcal{A} \cap \mathbb{S}_+^N$  satisfies (275) and (276). (confer §4.1.2.2) A matrix  $X \triangleq R^T R$  is an extreme point if and only if the smallest face, that contains  $X$ , of  $\mathcal{A} \cap \mathbb{S}_+^N$  has dimension 0; [268, §2.4] [269] *id est*, iff

$$\dim \mathcal{F}((\mathcal{A} \cap \mathbb{S}_+^N) \ni X) = \text{rank}(X)(\text{rank}(X) + 1)/2 - \text{rank}[\text{svec } RA_1R^T \text{ svec } RA_2R^T \dots \text{ svec } RA_mR^T] \quad (277)$$

(174) equals 0 in isomorphic  $\mathbb{R}^{N(N+1)/2}$ .

Now the intersection  $\mathcal{A} \cap \mathbb{S}_+^N$  is assumed bounded: Assume a given nonzero upper bound  $\rho$  on rank, a number of equalities

$$m = (\rho + 1)(\rho + 2)/2 \quad (278)$$

and matrix dimension  $N \geq \rho + 2 \geq 3$ . If the intersection is nonempty and bounded, then there exists a matrix  $X \in \mathcal{A} \cap \mathbb{S}_+^N$  such that

$$\text{rank } X \leq \rho \quad (279)$$

This represents a tightening of the upper bound; a reduction by exactly 1 of the bound provided by (275) given the same specified number  $m$  (278) of equalities; *id est*,

$$\text{rank } X \leq \frac{\sqrt{8m+1}-1}{2} - 1 \quad (280)$$

$\diamond$

---

<sup>2.53</sup> §4.1.2.2 contains an intuitive explanation. This bound is itself limited above, of course, by  $N$ ; a *tight* limit corresponding to an interior point of  $\mathbb{S}_+^N$ .

## 2.10 Conic independence (c.i.)

In contrast to extreme direction, the property *conically independent direction* is more generally applicable; inclusive of all closed convex cones (not only pointed closed convex cones). Arbitrary given directions  $\{\Gamma_i \in \mathbb{R}^n, i=1 \dots N\}$  comprise a *conically independent set* if and only if (*confer §2.1.2, §2.4.2.3*)

$$\Gamma_i \zeta_i + \dots + \Gamma_j \zeta_j - \Gamma_\ell = \mathbf{0}, \quad i \neq \dots \neq j \neq \ell = 1 \dots N \quad (281)$$

has no solution  $\zeta \in \mathbb{R}_+^N$  ( $\zeta_i \in \mathbb{R}_+$ ); in words, iff no direction from the given set can be expressed as a conic combination of those remaining; *e.g.*, Figure 52 [413, *conic independence test* (281) MATLAB]. Arranging any set of generators for a particular closed convex cone in a matrix columnar,

$$X \triangleq [\Gamma_1 \ \Gamma_2 \ \dots \ \Gamma_N] \in \mathbb{R}^{n \times N} \quad (282)$$

then this test of conic independence (281) may be expressed as a set of linear *feasibility problems*: for  $\ell = 1 \dots N$

$$\begin{array}{ll} \text{find} & \zeta \in \mathbb{R}^N \\ \text{subject to} & X\zeta = \Gamma_\ell \\ & \zeta \succeq 0 \\ & \zeta_\ell = 0 \end{array} \quad (283)$$

If feasible for any particular  $\ell$ , then the set is not conically independent.

To find all conically independent directions from a set via (283), generator  $\Gamma_\ell$  must be removed from the set once it is found (feasible) conically dependent on remaining generators in  $X$ . So, to continue testing remaining generators when  $\Gamma_\ell$  is found to be dependent,  $\Gamma_\ell$  must be discarded from matrix  $X$  before proceeding. A generator  $\Gamma_\ell$  that is instead found conically independent of remaining generators in  $X$ , on the other hand, is conically independent of any subset of remaining generators. A c.i. set thus found is not necessarily unique.

It is evident that linear independence (l.i.) of  $N$  directions implies their conic independence;

- l.i.  $\Rightarrow$  c.i.

which suggests, number of l.i. generators in the columns of  $X$  cannot exceed number of c.i. generators. Denoting by  $\mathbf{k}$  the number of conically independent generators contained in  $X$ , we have the most fundamental rank inequality for convex cones

$$\dim \text{aff } \mathcal{K} = \dim \text{aff}[\mathbf{0} \ X] = \text{rank } X \leq \mathbf{k} \leq N \quad (284)$$

Whereas  $N$  directions in  $n$  dimensions can no longer be linearly independent once  $N$  exceeds  $n$ , conic independence remains possible:

### 2.10.0.0.1 Table: Maximum number of c.i. directions

dimension $n$	sup $\mathbf{k}$ (pointed)	sup $\mathbf{k}$ (not pointed)
0	0	0
1	1	2
2	2	4
3	$\infty$	$\infty$
$\vdots$	$\vdots$	$\vdots$

Assuming veracity of this table, there is an apparent vastness between two and three dimensions. These numbers of conically independent directions indicate:

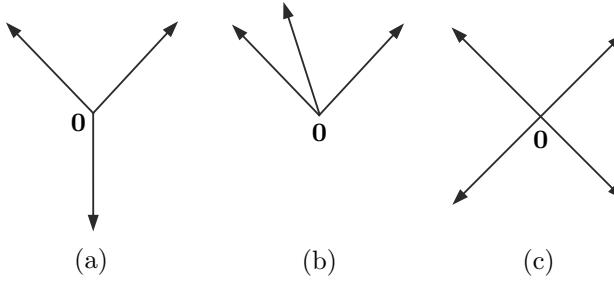


Figure 52: Vectors in  $\mathbb{R}^2$ : (a) affinely and conically independent, (b) affinely independent but not conically independent, (c) conically independent but not affinely independent. None of the examples exhibits linear independence. (In general, a.i.  $\nLeftrightarrow$  c.i.)

- Convex cones in dimensions 0, 1, and 2 must be polyhedral (§2.12.1).
- Full-dimensional pointed closed convex cones in dimensions 1 and 2 must be simplicial. (§2.8.1.1 p.87, §2.12.3.1.1)
- Pointed polyhedral cones in dimension 3 and higher can have an infinite number of faces; *id est*, can be full-dimensional and nonsimplicial (§2.12.3.1.1); *e.g.*, Figure 53.

It is also evident that dimension of Euclidean space cannot exceed number of conically independent directions possible;

$$\bullet \quad n \leq \sup k$$

Conic independence is certainly one convex idea that cannot be completely explained by a two-dimensional picture as Barvinok suggests (p.31) [27, p.vii].

### 2.10.1 Preservation of conic independence

Independence in the linear (§2.1.2.1), affine (§2.4.2.4), and conic senses can be preserved under linear transformation. Suppose a matrix  $X \in \mathbb{R}^{n \times N}$  (282) holds a conically independent set columnar. Consider a transformation on the domain of such matrices

$$T(X) : \mathbb{R}^{n \times N} \rightarrow \mathbb{R}^{n \times N} \triangleq XY \quad (285)$$

where fixed matrix  $Y \triangleq [y_1 \ y_2 \ \dots \ y_N] \in \mathbb{R}^{N \times N}$  represents linear operator  $T$ . Conic independence of  $\{Xy_i \in \mathbb{R}^n, i=1 \dots N\}$  demands, by definition (281),

$$Xy_i \zeta_i + \dots + Xy_j \zeta_j - Xy_\ell = \mathbf{0}, \quad i \neq \dots \neq j \neq \ell = 1 \dots N \quad (286)$$

have no solution  $\zeta \in \mathbb{R}_+^N$ . That is ensured by conic independence of  $\{y_i \in \mathbb{R}^N\}$  and by  $\mathcal{R}(Y) \cap \mathcal{N}(X) = \mathbf{0}$ ; seen by factoring out  $X$ .

#### 2.10.1.1 linear maps of cones

[22, §7] If  $\mathcal{K}$  is a convex cone in Euclidean space  $\mathcal{R}$  and  $T$  is any linear mapping from  $\mathcal{R}$  to Euclidean space  $\mathcal{M}$ , then  $T(\mathcal{K})$  is a convex cone in  $\mathcal{M}$  and  $x \preceq y$  with respect to  $\mathcal{K}$  implies  $T(x) \preceq T(y)$  with respect to  $T(\mathcal{K})$ . If  $\mathcal{K}$  is full-dimensional in  $\mathcal{R}$ , then so is  $T(\mathcal{K})$  in  $\mathcal{M}$ .

If  $T$  is a linear bijection, then  $x \preceq y \Leftrightarrow T(x) \preceq T(y)$ . If  $\mathcal{K}$  is pointed, then so is  $T(\mathcal{K})$ . And if  $\mathcal{K}$  is closed, so is  $T(\mathcal{K})$ . If  $\mathcal{F}$  is a face of  $\mathcal{K}$ , then  $T(\mathcal{F})$  is a face of  $T(\mathcal{K})$ .

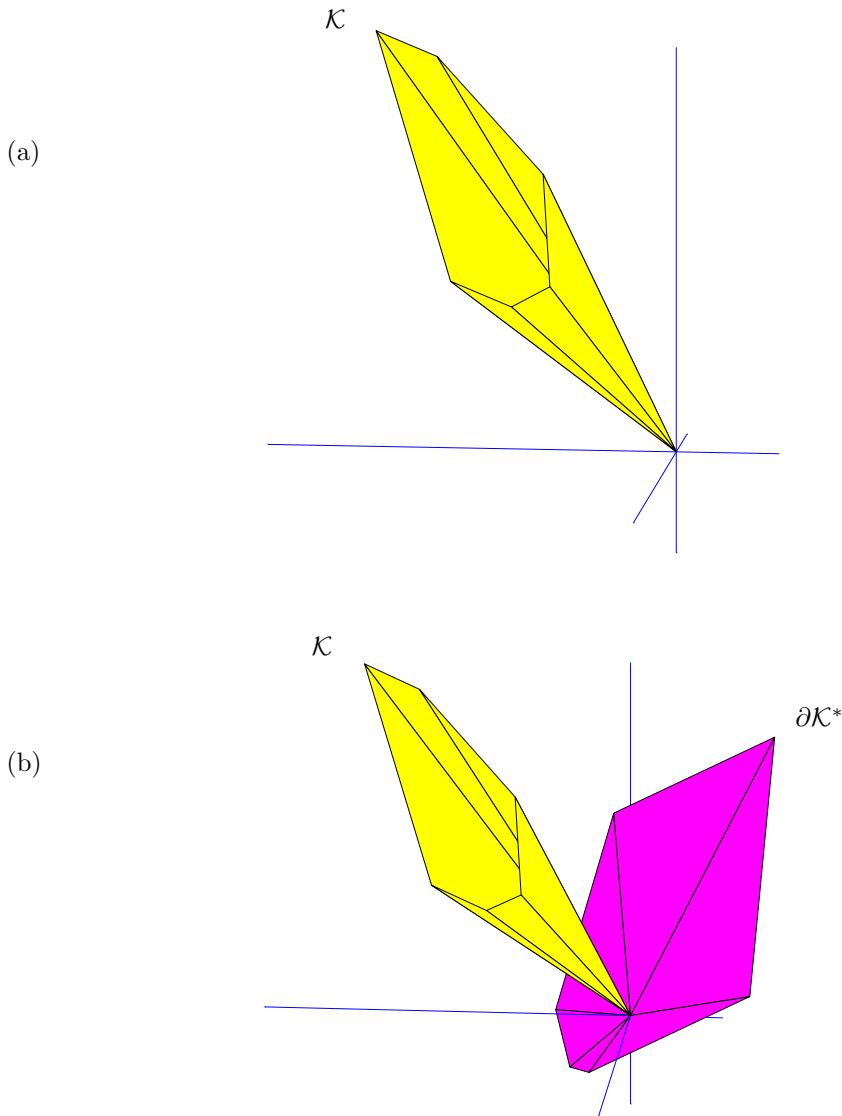


Figure 53: (a) A nonsimplicial pointed polyhedral cone (drawn truncated) in  $\mathbb{R}^3$  having six facets. The extreme directions, corresponding to six edges emanating from the origin, are generators for this cone; not linearly independent but they must be conically independent. (b) Boundary of dual cone  $\mathcal{K}^*$  (drawn truncated) is now added to the drawing of same  $\mathcal{K}$ .  $\mathcal{K}^*$  is polyhedral, proper, and has the same number of extreme directions as  $\mathcal{K}$  has facets.

Linear bijection is only a sufficient condition for pointedness and closedness; convex polyhedra (§2.12) are invariant to any linear or inverse linear transformation [27, §I.9] [343, p.44, thm.19.3].

### 2.10.2 Pointed closed convex $\mathcal{K}$ & conic independence

The following bullets can be derived from definitions (189) and (281) in conjunction with the *extremes theorem* (§2.8.1.1.1):

The set of all extreme directions from a pointed closed convex cone  $\mathcal{K} \subset \mathbb{R}^n$  is not necessarily a linearly independent set, yet it must be a conically independent set; (compare Figure 27 on page 59 with Figure 53a)

- $\{\text{extreme directions}\} \Rightarrow \{\text{c.i.}\}$

When a conically independent set of directions from pointed closed convex cone  $\mathcal{K}$  is known to comprise generators, conversely, then all directions from that set must be extreme directions of the cone;

- $\{\text{extreme directions}\} \Leftrightarrow \{\text{c.i. generators of pointed closed convex } \mathcal{K}\}$

Barker & Carlson [22, §1] call the extreme directions a *minimal generating set*<sup>2.54</sup> for a pointed closed convex cone. A minimal set of generators is therefore a conically independent set of generators, and *vice versa*,<sup>2.55</sup> for a pointed closed convex cone.

An arbitrary collection of  $n$  or fewer distinct extreme directions, from pointed closed convex cone  $\mathcal{K} \subset \mathbb{R}^n$ , is not necessarily a linearly independent set; *e.g.*, dual extreme directions (488) from Example 2.13.12.0.3.

- $\{\leq n \text{ extreme directions in } \mathbb{R}^n\} \not\Rightarrow \{\text{l.i.}\}$

Linear dependence of few extreme directions is another convex idea that cannot be explained by a two-dimensional picture as Barvinok suggests (p.31) [27, p.vii]; indeed, it only first comes to light in four dimensions! But there is a converse: [366, §2.10.9]

- $\{\text{extreme directions}\} \Leftarrow \{\text{l.i. generators of closed convex } \mathcal{K}\}$

#### 2.10.2.0.1 Example. Vertex-description of halfspace $\mathcal{H}$ about origin.

From  $n+1$  points in  $\mathbb{R}^n$  we can make a vertex-description of a convex cone that is a halfspace  $\mathcal{H}$ , where  $\{x_\ell \in \mathbb{R}^n, \ell=1 \dots n\}$  constitutes a minimal set of generators for a hyperplane  $\partial\mathcal{H}$  through the origin. An example is illustrated in Figure 54. By demanding the augmented set  $\{x_\ell \in \mathbb{R}^n, \ell=1 \dots n+1\}$  be affinely independent (we want vector  $x_{n+1}$  not parallel to  $\partial\mathcal{H}$ ), then

$$\begin{aligned} \mathcal{H} &= \bigcup_{\zeta \geq 0} (\zeta x_{n+1} + \partial\mathcal{H}) \\ &= \{\zeta x_{n+1} + \text{cone}\{x_\ell \in \mathbb{R}^n, \ell=1 \dots n\} \mid \zeta \geq 0\} \\ &= \text{cone}\{x_\ell \in \mathbb{R}^n, \ell=1 \dots n+1\} \end{aligned} \tag{287}$$

a union of parallel hyperplanes. Cardinality is one step beyond dimension of the ambient space, but  $\{x_\ell \forall \ell\}$  is a minimal set of generators for this convex cone  $\mathcal{H}$  which has no extreme elements.  $\square$

<sup>2.54</sup>A minimal generating set for any polyhedral cone (§2.12.1) is known as a *frame*; *e.g.*, Figure 54.

<sup>2.55</sup>This converse does not hold for nonpointed closed convex cones as Table 2.10.0.0.1 implies; *e.g.*, ponder four conically independent generators for a plane ( $n=2$ , Figure 52).

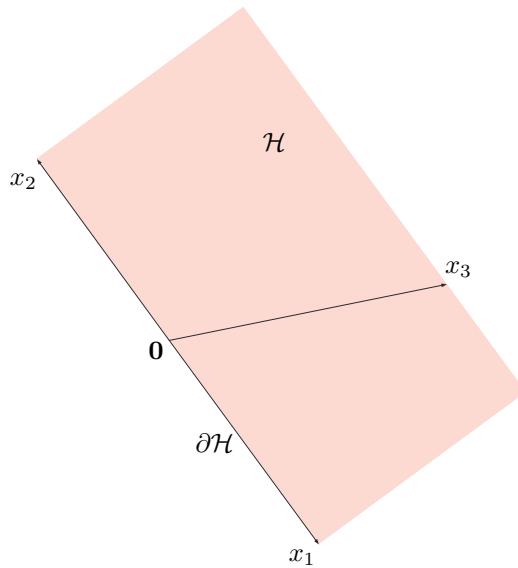


Figure 54: Minimal set of generators  $X = [x_1 \ x_2 \ x_3] \in \mathbb{R}^{2 \times 3}$  (not extreme directions) for halfspace about origin; affinely and conically independent. Any halfspace, about origin, is a polyhedral cone (§2.12.1) but is not pointed.

#### 2.10.2.0.2 Exercise. *Enumerating conically independent directions.*

Do Example 2.10.2.0.1 in  $\mathbb{R}$  and  $\mathbb{R}^3$  by drawing two figures corresponding to Figure 54 and enumerating  $n + 1$  conically independent generators for each. Describe a nonpointed polyhedral cone in three dimensions having more than eight conically independent generators. (confer Table 2.10.0.0.1) ▼

### 2.10.3 Utility of conic independence

Perhaps the most useful application of conic independence is determination of the intersection of closed convex cones from their halfspace-descriptions, or representation of the sum of closed convex cones from their vertex-descriptions.

- $\cap \mathcal{K}_i$  A halfspace-description for the intersection of any number of closed convex cones  $\mathcal{K}_i$  can be acquired by pruning normals; specifically, only the conically independent normals from the aggregate of all the halfspace-descriptions need be retained.
- $\sum \mathcal{K}_i$  Generators for the sum of any number of closed convex cones  $\mathcal{K}_i$  can be determined by retaining only the conically independent generators from the aggregate of all the vertex-descriptions.

Such conically independent sets are not necessarily unique or minimal.

## 2.11 When extreme means exposed

For any convex full-dimensional polyhedral set in  $\mathbb{R}^n$ , distinction between the terms *extreme* and *exposed* vanishes [366, §2.4] [126, §2.2] for faces of all dimensions except  $n$ ; their meanings become equivalent as we saw in Figure 22 (discussed in §2.6.1.2). In other words, each and every face of any polyhedral set (except the set itself) can be exposed by a hyperplane, and *vice versa*; e.g, Figure 27.

Lewis [276, §6] [242, §2.3.4] claims nonempty extreme proper subsets and the exposed subsets coincide for  $\mathbb{S}_+^n$ ; *id est*, each and every face of the positive semidefinite cone (whose dimension is less than dimension of the cone) is exposed. A more general discussion of cones having this property can be found in [379]; *e.g.*, Lorentz cone (181) [21, §II.A].

## 2.12 Convex polyhedra

Every polyhedron, such as the convex hull (87) of a bounded list  $X$ , can be expressed as the solution set of a finite system of linear equalities and inequalities, and *vice versa*. [126, §2.2]

### 2.12.0.0.1 Definition. Convex polyhedra, halfspace-description.

A convex polyhedron is the intersection of a finite number of halfspaces and hyperplanes;

$$\mathcal{P} = \{y \mid Ay \succeq b, Cy = d\} \subseteq \mathbb{R}^n \quad (288)$$

where coefficients  $A$  and  $C$  generally denote matrices. Each row of  $C$  is a vector normal to a hyperplane, while each row of  $A$  is a vector inward-normal to a hyperplane partially bounding a halfspace.  $\triangle$

By the *halfspaces theorem* in §2.4.1.1.1, a polyhedron thus described is a closed convex set possibly not full-dimensional; *e.g.*, Figure 22. Convex polyhedra<sup>2.56</sup> are finite-dimensional comprising all affine sets (§2.3.1, §2.1.4), polyhedral cones, line segments, rays, halfspaces, convex polygons, *solids* [250, def.104/6 p.343], polychora, *polytopes*,<sup>2.57</sup> *etcetera*.

It follows from definition (288) by exposure that each face of a convex polyhedron is a convex polyhedron.

Projection of any polyhedron on a subspace remains a polyhedron. More generally, image and inverse image of a convex polyhedron under any linear transformation remains a convex polyhedron; [27, §I.9] [343, thm.19.3] the foremost consequence being, invariance of polyhedral set closedness.

When  $b$  and  $d$  in (288) are  $\mathbf{0}$ , the resultant is a polyhedral cone. The set of all polyhedral cones is a subset of convex cones:

### 2.12.1 Polyhedral cone

From our study of cones, we see: the number of intersecting hyperplanes and halfspaces constituting a convex cone is possibly but not necessarily infinite. When the number is finite, the convex cone is termed *polyhedral*. That is the primary distinguishing feature between the set of all convex cones and polyhedra; all polyhedra, including polyhedral cones, are *finitely generated* [343, §19]. (Figure 55) We distinguish polyhedral cones in the set of all convex cones for this reason, although all convex cones of dimension 2 or less are polyhedral.

#### 2.12.1.0.1 Definition. Polyhedral cone, halfspace-description.<sup>2.58</sup> (confer(105))

A polyhedral cone is the intersection of a finite number of halfspaces and hyperplanes

<sup>2.56</sup>We consider only convex polyhedra throughout, but acknowledge the existence of concave polyhedra. [436, *Kepler-Poinsot Solid*]

<sup>2.57</sup>Some authors distinguish bounded polyhedra via designation *polytope*. [126, §2.2]

<sup>2.58</sup>Rockafellar [343, §19] proposes affine sets be handled via complementary pairs of affine inequalities; *e.g.*, antisymmetry  $Cy \succeq d$  and  $Cy \preceq d$  which can present severe difficulty to some *interior-point methods* of numerical solution.

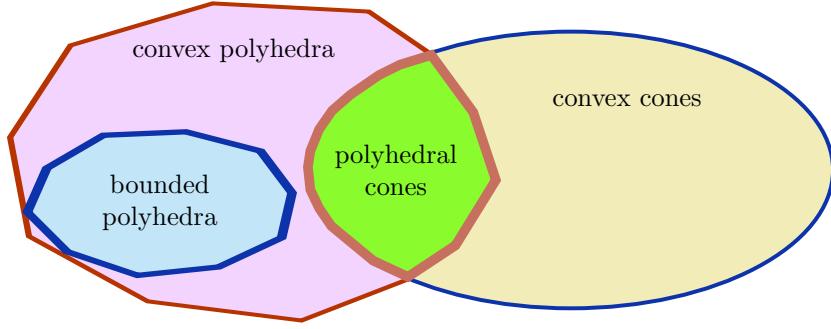


Figure 55: Polyhedral cones are finitely generated, unbounded, and convex.

about the origin;

$$\begin{aligned}
 \mathcal{K} &= \{y \mid Ay \succeq 0, Cy = \mathbf{0}\} \subseteq \mathbb{R}^n & (a) \\
 &= \{y \mid Ay \succeq 0, Cy \succeq 0, Cy \preceq 0\} & (b) \\
 &= \left\{ y \mid \begin{bmatrix} A \\ C \\ -C \end{bmatrix} y \succeq 0 \right\} & (c)
 \end{aligned} \tag{289}$$

where coefficients  $A$  and  $C$  generally denote matrices of finite dimension. Each row of  $C$  is a vector normal to a hyperplane containing the origin, while each row of  $A$  is a vector inward-normal to a hyperplane containing the origin and partially bounding a halfspace.  $\triangle$

A polyhedral cone thus defined is closed, convex (§2.4.1.1), has only a finite number of generators (§2.8.1.2), and can be not full-dimensional. (Minkowski) Conversely, any finitely generated convex cone is polyhedral. (Weyl) [366, §2.8] [343, thm.19.1]

#### 2.12.1.0.2 Exercise. *Unbounded convex polyhedra.*

Illustrate an unbounded polyhedron that is not a cone or its translation.  $\blacktriangledown$

From the definition it follows that any single hyperplane, through the origin, or any halfspace partially bounded by a hyperplane through the origin is a polyhedral cone. The most familiar example of polyhedral cone is any quadrant (or orthant, §2.1.3) generated by Cartesian half-axes. Esoteric examples of polyhedral cone include the point at the origin, any line through the origin, any ray having the origin as base such as the nonnegative real line  $\mathbb{R}_+$  in subspace  $\mathbb{R}$ , polyhedral flavor (proper) Lorentz cone (303), any subspace, and  $\mathbb{R}^n$ . More polyhedral cones are illustrated in Figure 53 and Figure 27.

## 2.12.2 Vertices of convex polyhedra

By definition, a vertex (§2.6.1.0.1) always lies on the relative boundary of a convex polyhedron. [250, def.115/6 p.358] In Figure 22, each vertex of the polyhedron is located at an intersection of three or more facets, and every edge belongs to precisely two facets [27, §VI.1 p.252]. In Figure 27, the only vertex of that polyhedral cone lies at the origin.

The set of all polyhedral cones is clearly a subset of convex polyhedra and a subset of convex cones (Figure 55). Not all convex polyhedra are bounded; evidently, neither can they all be described by the convex hull of a bounded set of points as defined in (87).

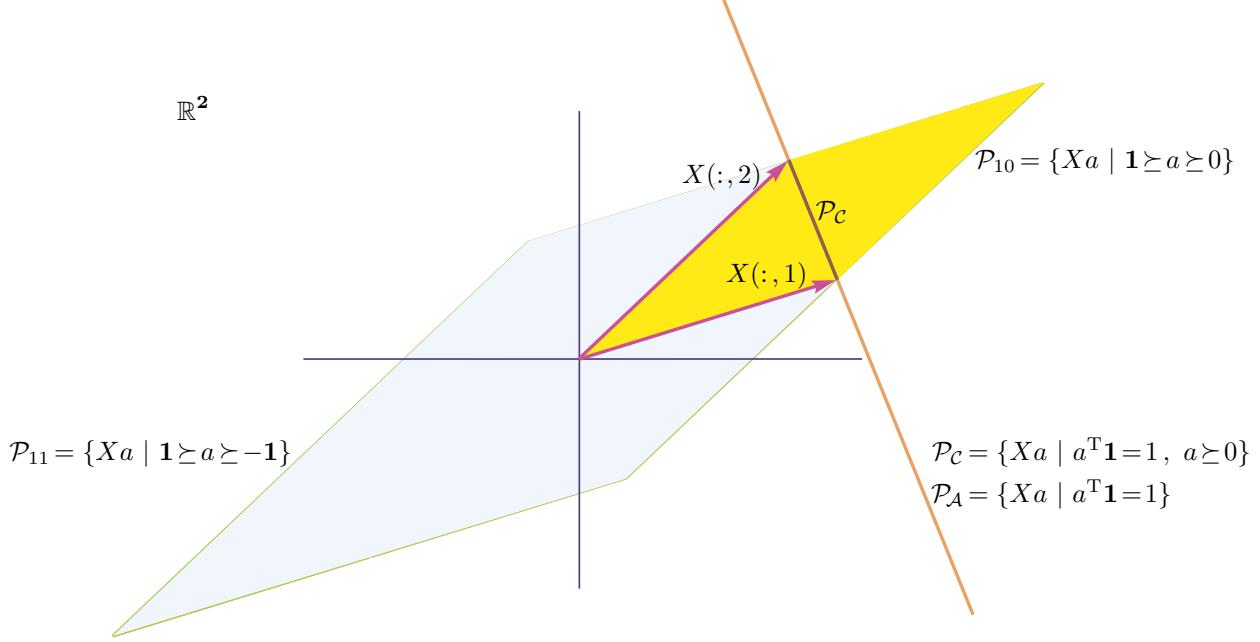


Figure 56: A polyhedron's generating list  $X$  does not necessarily constitute its vertices. Convex polyhedra  $\mathcal{P}_{11} \supset \mathcal{P}_{10}$  are *parallelograms*, polyhedron  $\mathcal{P}_A$  is a line, polyhedron  $\mathcal{P}_C$  is a line segment in  $\mathcal{P}_A \supset \mathcal{P}_C$ . In higher dimension,  $\mathcal{P}_{11}$  and  $\mathcal{P}_{10}$  are known as *parallelepipeds*. Were vector  $a$  unbounded above,  $\mathcal{P}_{10}$  would become a polyhedral cone  $\mathcal{K}$ . Were vector  $a$  unbounded above and below,  $\mathcal{P}_{11}$  and  $\mathcal{P}_{10}$  would become subspace  $\mathcal{R}(X)$ .

Hence a universal vertex-description of polyhedra in terms of that same finite-length list  $X$  (77):

#### 2.12.2.0.1 Definition. Convex polyhedra, vertex-description.

Denote upper  $u$  and lower  $\ell$  real vector bounds and truncated  $N$ -dimensional  $a$ -vector by

$$a_{i:j} = \begin{bmatrix} a_i \\ \vdots \\ a_j \end{bmatrix} \quad (290)$$

By discriminating a suitable finite-length *generating list* (or set) arranged columnar in  $X \in \mathbb{R}^{n \times N}$ , then any particular polyhedron may be described

$$\mathcal{P} = \{Xa \mid a_{1:k}^T \mathbf{1} = 1, u \succeq a_{m:N} \succeq \ell, \{1 \dots k\} \cup \{m \dots N\} = \{1 \dots N\}\} \quad (291)$$

where  $0 \leq k \leq N$  and  $1 \leq m \leq N+1$ . Setting  $k=0$  removes the affine equality condition. Setting  $m=N+1$  removes the inequality.  $\triangle$

Coefficient indices in (291) may or may not be overlapping. From (79), (87), (105), and (144), we summarize how the coefficient conditions may be applied;

subspace	$\longrightarrow$	$\infty \succeq a \succeq -\infty$	(292)
parallelepiped	$\longrightarrow$	$u \succeq a \succeq \ell$	
affine set	$\longrightarrow$	$a_{1:k}^T \mathbf{1} = 1$	
polyhedral cone	$\longrightarrow$	$a_{m:N} \succeq 0$	

It is always possible to describe a convex hull in a region of overlapping indices because, for  $1 \leq m \leq k \leq N$

$$\{a_{m:k} \mid a_{m:k}^T \mathbf{1} = 1, a_{m:k} \succeq 0\} \subseteq \{a_{m:k} \mid a_{1:k}^T \mathbf{1} = 1, a_{m:N} \succeq 0\} \quad (293)$$

Generating list members are neither unique or necessarily vertices of the corresponding polyhedron; *e.g.*, Figure 56. Indeed, for convex hull (87) (a special case of (291)), some subset of list members may reside in the polyhedron's relative interior. Conversely, convex hull of the vertices and extreme rays of a polyhedron is identical to the convex hull of any list generating that polyhedron; that is, *extremes theorem* 2.8.1.1.1.

### 2.12.2.1 Vertex-description of polyhedral cone

Given closed convex cone  $\mathcal{K}$  in a subspace of  $\mathbb{R}^n$  having any set of generators for it arranged in a matrix  $X \in \mathbb{R}^{n \times N}$  as in (282), then that cone is described setting  $m=1$  and  $k=0$  in vertex-description (291): (*confer*(289))

$$\mathcal{K} = \text{cone } X = \{Xa \mid a \succeq 0\} \subseteq \mathbb{R}^n \quad (105)$$

a conic hull of  $N$  generators.

### 2.12.2.2 Pointedness

(§2.7.2.1.2) [366, §2.10] Assuming all generators constituting the columns of  $X \in \mathbb{R}^{n \times N}$  are nonzero, polyhedral cone  $\mathcal{K}$  is pointed if and only if there is no nonzero  $a \succeq 0$  that solves  $Xa = \mathbf{0}$ , or iff<sup>2.59</sup>

$$\mathcal{N}(X) \cap \mathbb{R}_+^N = \mathbf{0} \quad (294)$$

or iff

$$\begin{array}{ll} \text{find} & a \\ \text{subject to} & \begin{array}{l} Xa = \mathbf{0} \\ \mathbf{1}^T a = 1 \\ a \succeq 0 \end{array} \end{array} \quad (295)$$

is infeasible. Otherwise, the cone will contain at least one line and there can be no vertex nor extreme direction; *id est*, the cone cannot, otherwise, be pointed. Any subspace, Euclidean vector space  $\mathbb{R}^n$ , or any halfspace are examples of nonpointed polyhedral cone.

This null-pointedness criterion  $Xa = \mathbf{0}$  means that a pointed polyhedral cone is invariant to linear injective transformation. Examples of pointed polyhedral cone  $\mathcal{K}$  include: the origin, any  $\mathbf{0}$ -based ray in a subspace, any two-dimensional V-shaped cone in a subspace, any orthant in  $\mathbb{R}^n$  or  $\mathbb{R}^{m \times n}$ ; *e.g.*, nonnegative real line  $\mathbb{R}_+$  in vector space  $\mathbb{R}$ .

### 2.12.3 Unit simplex

A peculiar subset of the nonnegative orthant with halfspace-description

$$\mathcal{S} \triangleq \{s \mid s \succeq 0, \mathbf{1}^T s \leq 1\} \subseteq \mathbb{R}_+^n \quad (296)$$

is a unique bounded convex full-dimensional polyhedron called *unit simplex* (Figure 57) having  $n+1$  facets,  $n+1$  vertices, and dimension

$$\dim \mathcal{S} = n \quad (297)$$

---

<sup>2.59</sup>If  $\text{rank } X = n$ , then dual cone  $\mathcal{K}^*$  (§2.13.1) is pointed. (314) The intersection with  $\mathbb{R}_+^N$  means that an  $a$  vector in  $\mathcal{N}(X)$  can have nonnegative entries, just not exclusively (excepting the origin).

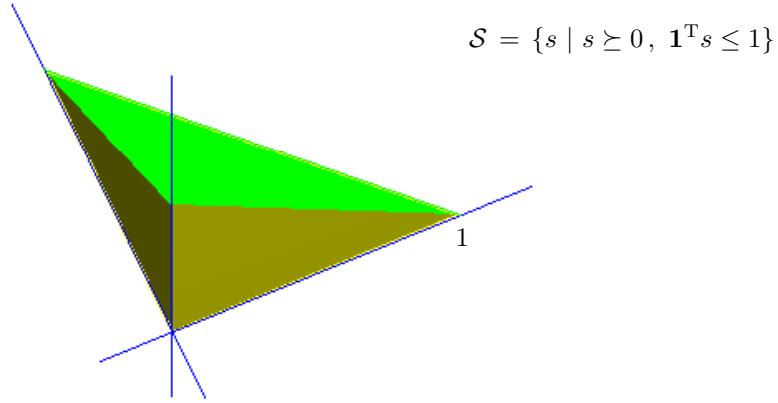


Figure 57: Unit simplex  $\mathcal{S}$  in  $\mathbb{R}^3$  is a unique solid tetrahedron but is not *regular*.

The origin supplies one vertex while heads of the *standard basis* [228] [368]  $\{e_i, i=1 \dots n\}$  in  $\mathbb{R}^n$  constitute those remaining; thus its vertex-description:

$$\begin{aligned}\mathcal{S} &= \text{conv} \{\mathbf{0}, \{e_i, i=1 \dots n\}\} \\ &= \{[\mathbf{0} \ e_1 \ e_2 \ \dots \ e_n] a \mid a^T \mathbf{1} = 1, a \succeq 0\}\end{aligned}\tag{298}$$

In  $\mathbb{R}^0$  the unit simplex is the point at the origin, in  $\mathbb{R}$  the unit simplex is the line segment  $[0, 1]$ , in  $\mathbb{R}^2$  it is a triangle and its relative interior, in  $\mathbb{R}^3$  it is the convex hull of a tetrahedron (Figure 57), in  $\mathbb{R}^4$  it is the convex hull of a *pentatope* [436], and so on.

### 2.12.3.1 Simplex

The unit simplex comes from a class of general polyhedra called *simplex*, having vertex-description: [103] [343] [434] [126] given  $n \geq k$

$$\text{conv}\{x_\ell \in \mathbb{R}^n \mid \ell = 1 \dots k+1, \dim \text{aff}\{x_\ell\} = k\}\tag{299}$$

So defined, a simplex is a closed bounded convex set possibly not full-dimensional. Examples of simplices, by increasing affine dimension, are: a point, any line segment, any triangle and its relative interior, a general tetrahedron, any five-vertex polychoron, and so on.

#### 2.12.3.1.1 Definition. Simplicial cone.

A proper (§2.7.2.2.1) polyhedral cone  $\mathcal{K}$  in  $\mathbb{R}^n$  is called *simplicial* iff  $\mathcal{K}$  has exactly  $n$  extreme directions; [21, §II.A] equivalently, iff pointed polyhedral cone  $\mathcal{K}$  in  $\mathbb{R}^n$  can be generated by  $n$  linearly independent directions.  $\triangle$

- simplicial cone  $\Rightarrow$  proper polyhedral cone

Whereas a ray having base  $\mathbf{0}$  in  $\mathbb{R}$  is a simplicial cone, any full-dimensional pointed closed convex cone in  $\mathbb{R}^2$  is simplicial. There are an infinite variety of simplicial cones in  $\mathbb{R}^n$ ; *e.g.*, Figure 27, Figure 58, Figure 68. Any orthant is simplicial, as is any rotation thereof.

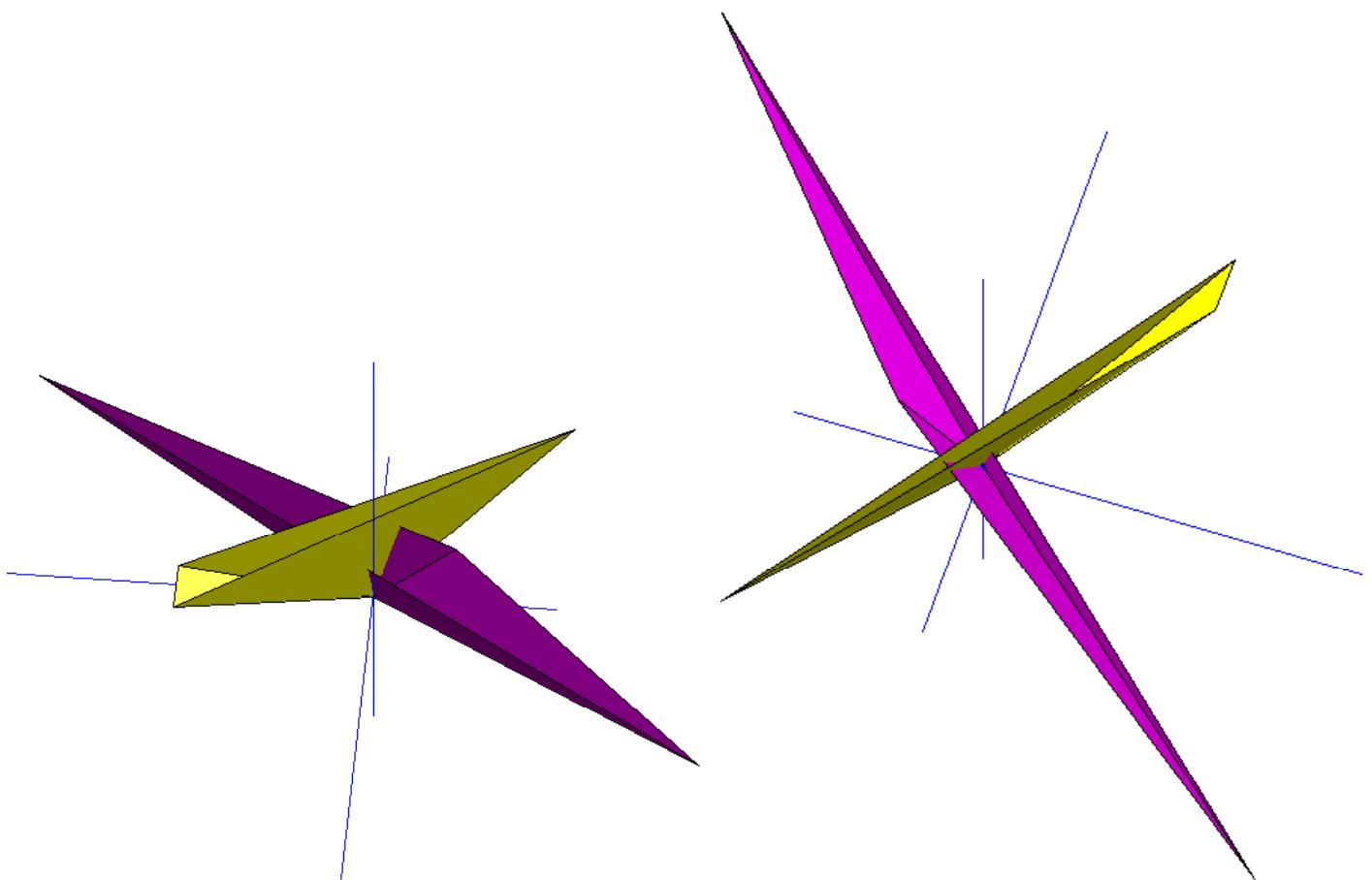


Figure 58: Two views of a simplicial cone and its dual in  $\mathbb{R}^3$ . Semiinfinite boundary of each cone is truncated for illustration. Each cone has three facets (*confer* §2.13.12.0.3). (Cartesian axes drawn for reference.)

### 2.12.4 Converting between descriptions

Conversion between halfspace- (288) (289) and vertex-description (87) (291) is nontrivial, in general, [16] [126, §2.2] [237] but the conversion is easy for simplices. [65, §2.2.4] Nonetheless, we tacitly assume the two descriptions of polyhedra to be equivalent. [343, §19 thm.19.1] We explore conversions in §2.13.4, §2.13.10, and §2.13.12:

## 2.13 Dual cone & generalized inequality & biorthogonal expansion

These three concepts, dual cone, generalized inequality, and biorthogonal expansion, are inextricably melded; meaning, it is difficult to completely discuss one without mentioning the others. The dual cone is critical in tests for convergence by contemporary primal/dual methods for numerical solution of conic problems. [453] [309, §4.5] For unique minimum-distance projection on a closed convex cone  $\mathcal{K}$ , the negative dual cone  $-\mathcal{K}^*$  plays the role that orthogonal complement plays for subspace projection.<sup>2.60</sup> (§E.9.2, Figure 193) Indeed,  $-\mathcal{K}^*$  is the algebraic complement in  $\mathbb{R}^n$ :

$$\mathcal{K} \boxplus -\mathcal{K}^* = \mathbb{R}^n \quad (2224)$$

where  $\boxplus$  denotes unique orthogonal vector sum.

One way to think of a pointed closed convex cone is as a new kind of coordinate system whose basis is generally nonorthogonal; a conic system, very much like the familiar Cartesian system whose analogous cone is the first quadrant (the nonnegative orthant). Generalized inequality  $\succeq_{\mathcal{K}}$  is a formalized means to determine membership to any pointed closed convex cone  $\mathcal{K}$  (§2.7.2.2) whereas *biorthogonal expansion* is, fundamentally, an expression of coordinates in a pointed conic system whose axes are linearly independent but not necessarily orthogonal. When cone  $\mathcal{K}$  is the nonnegative orthant, then these three concepts come into alignment with the Cartesian prototype: biorthogonal expansion becomes orthogonal expansion, the dual cone becomes identical to the orthant, and generalized inequality obeys a total order entrywise.

### 2.13.1 Dual cone

For any set  $\mathcal{K}$  (convex or not), its *dual cone* [122, §4.2]

$$\mathcal{K}^* \triangleq \{y \in \mathbb{R}^n \mid \langle y, x \rangle \geq 0 \text{ for all } x \in \mathcal{K}\} \quad (300)$$

is a unique cone<sup>2.61</sup> that is always closed and convex because it is an intersection of halfspaces (§2.4.1.1). Each halfspace has inward-normal  $x$ , belonging to  $\mathcal{K}$ , and boundary containing the origin; *e.g.*, Figure 59a.

When cone  $\mathcal{K}$  is convex, there is a second and equivalent construction: Dual cone  $\mathcal{K}^*$  is the union of each and every vector  $y$  inward-normal to a hyperplane supporting  $\mathcal{K}$  (§2.4.2.6.1); *e.g.*, Figure 59b. When  $\mathcal{K}$  is represented by a halfspace-description such as (289), for example, where

$$(150) \quad A \triangleq \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad C \triangleq \begin{bmatrix} c_1^T \\ \vdots \\ c_p^T \end{bmatrix} \in \mathbb{R}^{p \times n} \quad (301)$$

<sup>2.60</sup>Namely, projection on a subspace is ascertainable from projection on its orthogonal complement (Figure 192).

<sup>2.61</sup>The dual cone is the negative polar cone defined by many authors;  $\mathcal{K}^* = -\mathcal{K}^\circ$ . [225, §A.3.2] [343, §14] [42] [27] [366, §2.7]

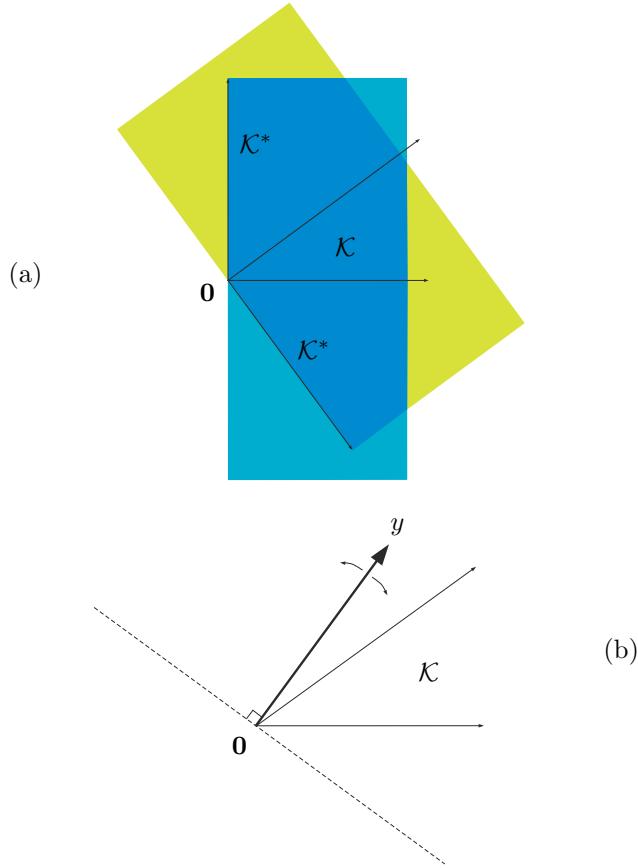


Figure 59: Two equivalent constructions of dual cone  $\mathcal{K}^*$  in  $\mathbb{R}^2$ : (a) Showing construction by intersection of halfspaces about  $\mathbf{0}$  (drawn truncated). Only those two halfspaces, whose bounding hyperplanes have inward-normal corresponding to an extreme direction of this pointed closed convex cone  $\mathcal{K} \subset \mathbb{R}^2$ , need be drawn; by (372). (b) Suggesting construction by union of inward-normals  $y$  to each and every hyperplane  $\underline{\partial\mathcal{H}_+}$  supporting  $\mathcal{K}$ . This interpretation is valid when  $\mathcal{K}$  is convex because existence of a supporting hyperplane is then guaranteed (§2.4.2.6).

then the dual cone can be represented as the conic hull

$$\mathcal{K}^* = \text{cone}\{a_1, \dots, a_m, \pm c_1, \dots, \pm c_p\} \quad (302)$$

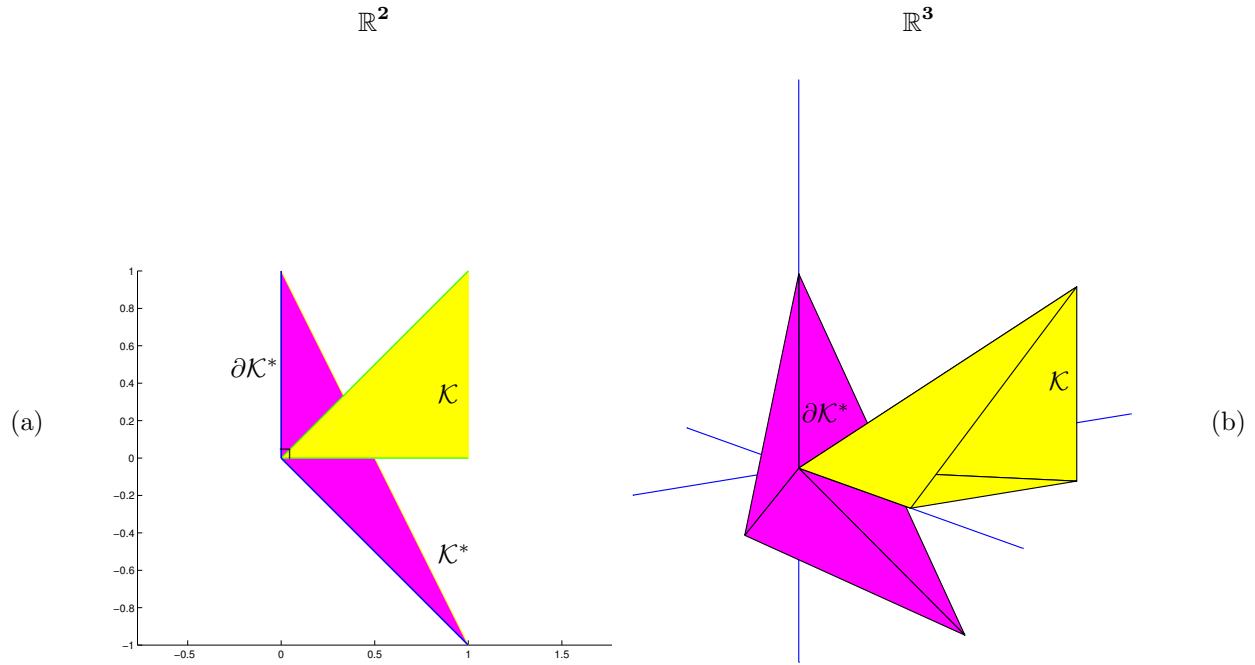
a vertex-description, because each and every conic combination of normals from the halfspace-description of  $\mathcal{K}$  yields another inward-normal to a hyperplane supporting  $\mathcal{K}$ .

$\mathcal{K}^*$  can also be constructed pointwise using projection theory from §E.9.2: for  $P_{\mathcal{K}}x$  the Euclidean projection of point  $x$  on closed convex cone  $\mathcal{K}$

$$-\mathcal{K}^* = \{x - P_{\mathcal{K}}x \mid x \in \mathbb{R}^n\} = \{x \in \mathbb{R}^n \mid P_{\mathcal{K}}x = \mathbf{0}\} \quad (2225)$$

#### 2.13.1.0.1 Exercise. Manual dual cone construction.

Perhaps the most instructive graphical method of dual cone construction is cut-and-try. Find the dual of each polyhedral cone from Figure 60 by using dual cone equation (300). ▼



$$x \in \mathcal{K} \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } y \in \mathcal{G}(\mathcal{K}^*) \quad (369)$$

Figure 60: Dual cone construction by right angle. Each extreme direction of a proper polyhedral cone is orthogonal to a facet of its dual cone, and *vice versa*, in any dimension. (§2.13.7.1) (a) This characteristic guides graphical construction of dual cone in two dimensions: It suggests finding dual-cone boundary  $\partial$  by making right angles with extreme directions of polyhedral cone. The construction is then pruned so that each dual boundary vector does not exceed  $\pi/2$  radians in angle with each and every vector from polyhedral cone. Were dual cone in  $\mathbb{R}^2$  to narrow, Figure 63 would be reached in limit. (b) Same polyhedral cone and its dual continued into three dimensions. (confer Figure 68)

### 2.13.1.0.2 Exercise. Dual cone definitions.

What is  $\{x \in \mathbb{R}^n \mid x^T z \geq 0 \quad \forall z \in \mathbb{R}^n\}$ ?

What is  $\{x \in \mathbb{R}^n \mid x^T z \geq 1 \quad \forall z \in \mathbb{R}^n\}$ ?

What is  $\{x \in \mathbb{R}^n \mid x^T z \geq 1 \quad \forall z \in \mathbb{R}_+^n\}$ ?



As defined, dual cone  $\mathcal{K}^*$  exists even when affine hull of the original cone is a proper subspace; *id est*, even when the original cone is not full-dimensional.<sup>2.62</sup>

#### 2.13.1.1 Examples of dual cone

When  $\mathcal{K}$  is  $\mathbb{R}^n$ ,  $\mathcal{K}^*$  is the point at the origin, and *vice versa*.

When  $\mathcal{K}$  is a subspace,  $\mathcal{K}^*$  is its orthogonal complement, and *vice versa*. (§E.9.2.1, Figure 61)

When cone  $\mathcal{K}$  is a halfspace in  $\mathbb{R}^n$  with  $n > 0$  (Figure 63 for example), the dual cone  $\mathcal{K}^*$  is a ray (base **0**) belonging to that halfspace but orthogonal to its bounding hyperplane (that contains the origin), and *vice versa*.

When convex cone  $\mathcal{K}$  is a closed halfplane in  $\mathbb{R}^3$  (Figure 62), it is neither pointed or full-dimensional; hence, the dual cone  $\mathcal{K}^*$  can be neither full-dimensional or pointed.

When  $\mathcal{K}$  is any particular orthant in  $\mathbb{R}^n$ , the dual cone is identical; *id est*,  $\mathcal{K} = \mathcal{K}^*$ .

When  $\mathcal{K}$  is any quadrant in subspace  $\mathbb{R}^2$ ,  $\mathcal{K}^*$  is a wedge-shaped polyhedral cone in  $\mathbb{R}^3$ ; e.g., for  $\mathcal{K}$  equal to quadrant I,  $\mathcal{K}^* = \begin{bmatrix} \mathbb{R}_+^2 \\ \mathbb{R} \end{bmatrix}$ .

When  $\mathcal{K}$  is a polyhedral flavor Lorentz cone (*confer* (181))

$$\mathcal{K}_\ell = \left\{ \begin{bmatrix} x \\ t \end{bmatrix} \in \mathbb{R}^n \times \mathbb{R} \mid \|x\|_\ell \leq t \right\}, \quad \ell \in \{1, \infty\} \quad (303)$$

its dual is the proper cone [65, exmp.2.25]

$$\mathcal{K}_q = \mathcal{K}_\ell^* = \left\{ \begin{bmatrix} x \\ t \end{bmatrix} \in \mathbb{R}^n \times \mathbb{R} \mid \|x\|_q \leq t \right\}, \quad \ell \in \{1, 2, \infty\} \quad (304)$$

where  $\|x\|_\ell^* = \|x\|_q$  is that norm dual to  $\|x\|_\ell$  determined via solution to  $1/\ell + 1/q = 1$ .<sup>2.63</sup> Figure 66 illustrates  $\mathcal{K} = \mathcal{K}_1$  and  $\mathcal{K}^* = \mathcal{K}_1^* = \mathcal{K}_\infty$  in  $\mathbb{R}^2 \times \mathbb{R}$ .

To further motivate our understanding of the dual cone, consider the ease with which convergence can be ascertained in the following optimization problem (306p):

#### 2.13.1.1.1 Example. Dual problem. *(confer* §4.1)

*Duality is a powerful and widely employed tool in applied mathematics for a number of reasons. First, the dual program is always convex even if the primal is not. Second, the number of variables in the dual is equal to the number of constraints in the primal which is often less than the number of variables in the primal program. Third, the maximum value achieved by the dual problem is often equal to the minimum of the primal.* – [335, §2.1.3] When not equal, the dual always provides a bound on the primal optimal objective. For convex problems, the dual variables provide necessary and sufficient optimality conditions:

Essentially, *Lagrange duality theory* concerns representation of a given optimization problem as half of a *minimax problem*. [343, §36] [65, §5.4] Given any real function  $f(x, z)$

$$\underset{x}{\text{minimize}} \underset{z}{\text{maximize}} f(x, z) \geq \underset{z}{\text{maximize}} \underset{x}{\text{minimize}} f(x, z) \quad (305)$$

<sup>2.62</sup>Rockafellar formulates dimension of  $\mathcal{K}$  and  $\mathcal{K}^*$ . [343, §14.6.1] His monumental work *Convex Analysis* has not one figure or illustration. See [27, §II.16] for illustration of Rockafellar's *recession cone* [43].

<sup>2.63</sup>Dual norm is not a conjugate or dual function.

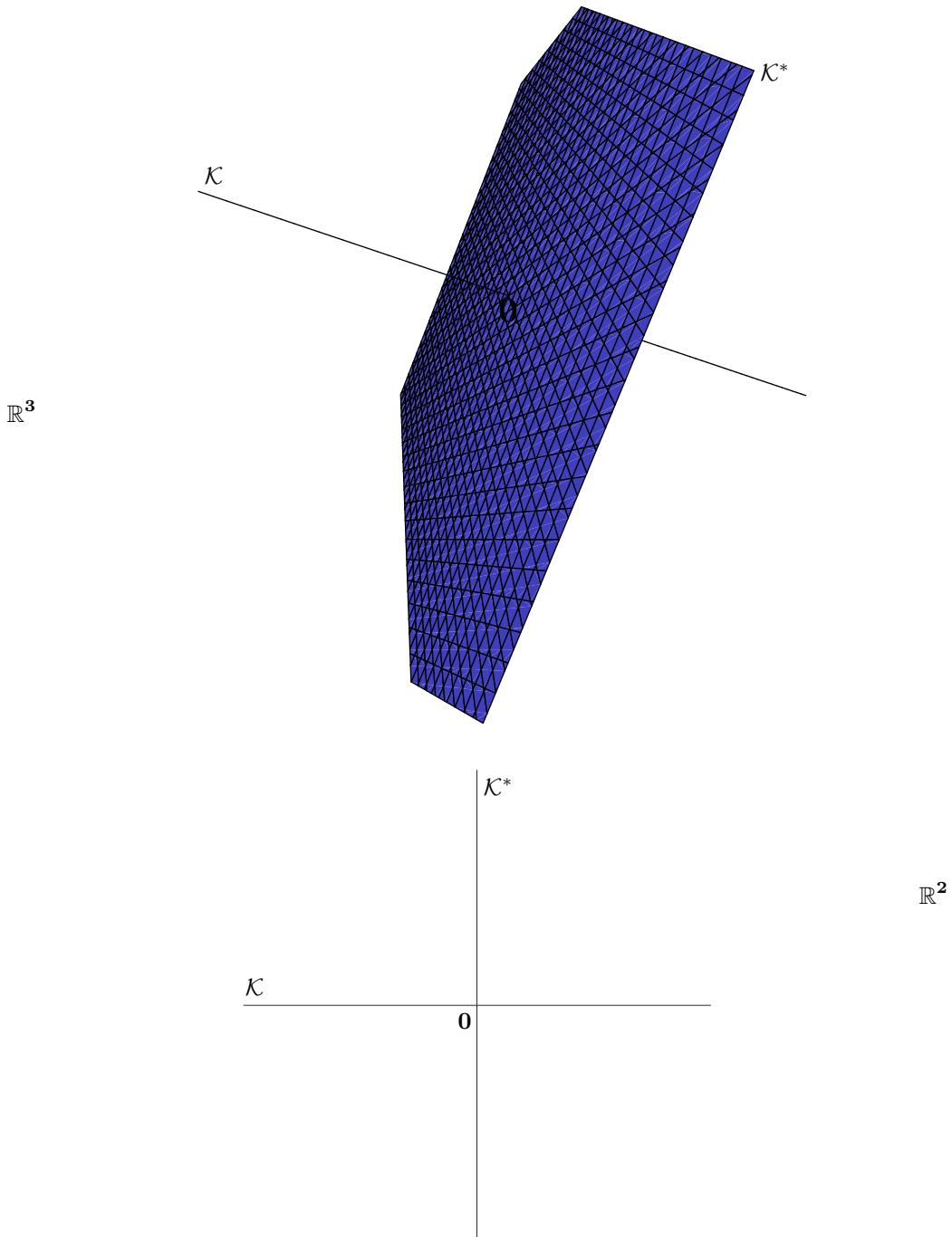


Figure 61: When convex cone  $\mathcal{K}$  is any one Cartesian axis, its dual cone is the convex hull of all axes remaining; its orthogonal complement. In  $\mathbb{R}^3$ , dual cone  $\mathcal{K}^*$  (drawn tiled and truncated) is a hyperplane through origin; its normal belongs to line  $\mathcal{K}$ . In  $\mathbb{R}^2$ , dual cone  $\mathcal{K}^*$  is a line through origin while convex cone  $\mathcal{K}$  is that line orthogonal.

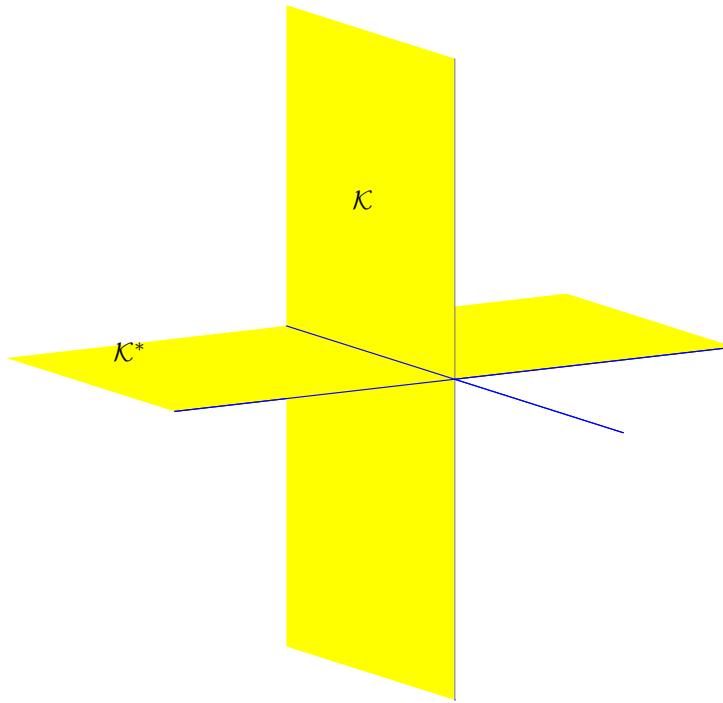


Figure 62:  $\mathcal{K}$  and  $\mathcal{K}^*$  are halfplanes in  $\mathbb{R}^3$ ; *blades*. Both semiinfinite convex cones appear truncated. Each cone is like  $\mathcal{K}$  from Figure 63 but embedded in a two-dimensional subspace of  $\mathbb{R}^3$ . (Cartesian coordinate axes drawn for reference.)

always holds. But when met with equality, then we have *strong duality* and a *saddle value* [173] exists. (Figure 64) [340, p.3]

Consider primal conic problem (p) (over cone  $\mathcal{K}$ ) and its corresponding *dual problem* (d): [327, §3.3.1] [268, §2.1] [269] given vectors  $\alpha, \beta$  and matrix constant  $C$

$$(p) \quad \begin{array}{ll} \text{minimize}_{\substack{x}} & \alpha^T x \\ \text{subject to} & x \in \mathcal{K} \\ & Cx = \beta \end{array} \quad (d) \quad \begin{array}{ll} \text{maximize}_{\substack{y, z}} & \beta^T z \\ \text{subject to} & y \in \mathcal{K}^* \\ & C^T z + y = \alpha \end{array}$$

Observe: the dual problem is also conic, and its objective function value never exceeds that of the primal;

$$\begin{aligned} \alpha^T x &\geq \beta^T z \\ x^T (C^T z + y) &\geq (Cx)^T z \\ x^T y &\geq 0 \end{aligned} \quad (307)$$

known as *weak duality* which holds by definition (300). Under the sufficient condition that (306p) is a *convex problem*<sup>2.64</sup> satisfying *Slater's condition* (p.229), then equality

$$x^* y^* = 0 \quad (308)$$

is achieved; which is necessary and sufficient for optimality (§2.13.11.1.5); each problem (p) and (d) attains the same optimal value of its objective, and each problem is called a

---

<sup>2.64</sup>In this context, problems (p) and (d) are convex if  $\mathcal{K}$  is a convex cone.

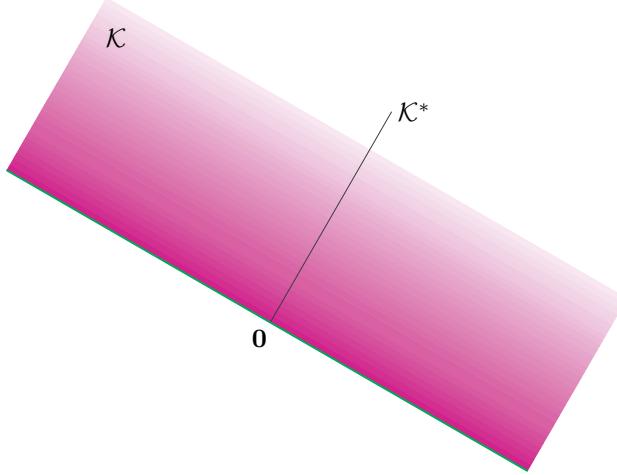


Figure 63: Polyhedral cone  $\mathcal{K}$  is a halfspace about origin in  $\mathbb{R}^2$ . Pointed dual cone  $\mathcal{K}^*$  is a ray base  $\mathbf{0}$ , hence not full-dimensional in  $\mathbb{R}^2$ ; so  $\mathcal{K}$  cannot be pointed, hence has no extreme directions nor vertex. (Both convex cones drawn truncated.)

*strong dual* to the other because the *duality gap* (optimal primal–dual objective difference) becomes 0. Then (p) and (d) are together equivalent to the minimax problem

$$\begin{array}{ll} \text{minimize}_{x,y,z} & \alpha^T x - \beta^T z \\ \text{subject to} & x \in \mathcal{K}, \quad y \in \mathcal{K}^* \\ & Cx = \beta, \quad C^T z + y = \alpha \end{array} \quad (\text{p}) - (\text{d}) \quad (309)$$

whose optimal objective always has the saddle value 0 (regardless of the particular convex cone  $\mathcal{K}$  and other problem parameters). [400, §3.2] Thus determination of convergence for either primal or dual problem is facilitated.

Were convex cone  $\mathcal{K}$  polyhedral (§2.12.1), then primal problem (p) and its dual (d) would be linear programs (LP). Selfdual nonnegative orthant  $\mathcal{K}$  yields the *prototypical primal linear program* and its dual. Were  $\mathcal{K}$  a positive semidefinite cone, then problem (p) has the form of *prototypical primal semidefinite program* (SDP (697)) with (d) its dual.

The dual problem may be solvable more quickly by computer. It is sometimes possible to solve a primal problem by way of its dual; advantageous *when the dual problem is easier to solve than the primal problem, for example, because it can be solved analytically, or has some special structure that can be exploited*. –[65, §5.5.5] (§4.2.3.1)  $\square$

### 2.13.1.2 Key properties of dual cone

- For any cone,  $(-\mathcal{K})^* = -\mathcal{K}^*$
- For any cones  $\mathcal{K}_1$  and  $\mathcal{K}_2$ ,  $\mathcal{K}_1 \subseteq \mathcal{K}_2 \Rightarrow \mathcal{K}_1^* \supseteq \mathcal{K}_2^*$  [366, §2.7]
- (Cartesian product) For closed convex cones  $\mathcal{K}_1$  and  $\mathcal{K}_2$ , their Cartesian product  $\mathcal{K} = \mathcal{K}_1 \times \mathcal{K}_2$  is a closed convex cone, and

$$\mathcal{K}^* = (\mathcal{K}_1 \times \mathcal{K}_2)^* = \mathcal{K}_1^* \times \mathcal{K}_2^* \quad (310)$$

where each dual is determined with respect to a cone's ambient space.

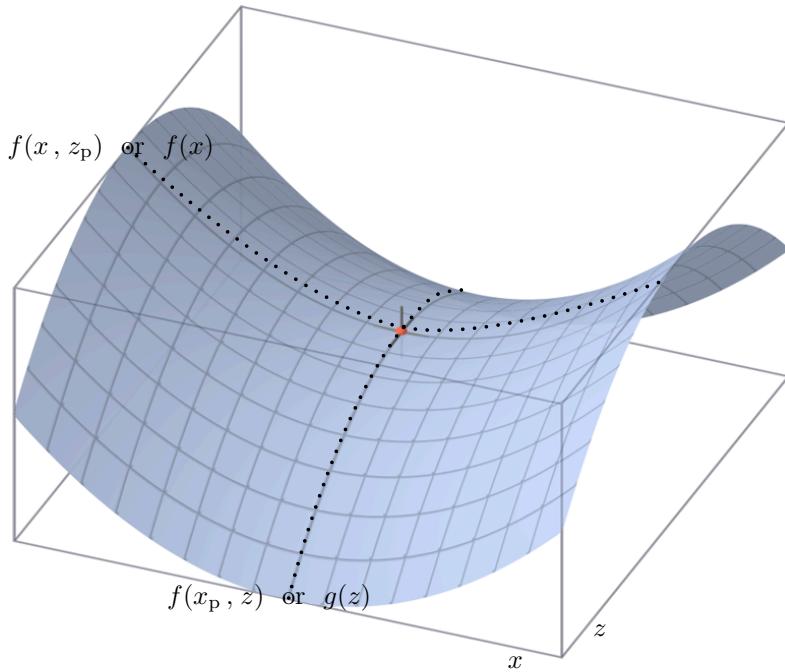


Figure 64: (Drawing by Lucas V. Barbosa.) This serves as mnemonic icon for primal and dual problems, although objective functions from conic problems (306p) (306d) are linear. When problems are *strong duals*, duality gap is 0; meaning, functions  $f(x)$ ,  $g(z)$  (dotted) *kiss* at saddle value as depicted at center. Otherwise, dual functions never meet ( $f(x) > g(z)$ ) by (305).

- (conjugation) [343, §14] [122, §4.5] [366, p.52] When  $\mathcal{K}$  is any convex cone, dual of the dual cone equals closure of the original cone;

$$\mathcal{K}^{**} = \overline{\mathcal{K}} \quad (311)$$

is the intersection of all halfspaces about the origin that contain  $\mathcal{K}$ . Because  $\mathcal{K}^{***} = \mathcal{K}^*$  always holds,

$$\mathcal{K}^* = (\overline{\mathcal{K}})^* \quad (312)$$

When convex cone  $\mathcal{K}$  is closed, then dual of the dual cone is the original cone;  $\mathcal{K}^{**} = \mathcal{K} \Leftrightarrow \mathcal{K}$  is a closed convex cone: [366, p.53, p.95]

$$\mathcal{K} = \{x \in \mathbb{R}^n \mid \langle y, x \rangle \geq 0 \ \forall y \in \mathcal{K}^*\} \quad (313)$$

- If any cone  $\mathcal{K}$  is full-dimensional, then  $\mathcal{K}^*$  is pointed;

$$\mathcal{K} \text{ full-dimensional} \Rightarrow \mathcal{K}^* \text{ pointed} \quad (314)$$

If the closure of any convex cone  $\mathcal{K}$  is pointed, conversely, then  $\mathcal{K}^*$  is full-dimensional;

$$\overline{\mathcal{K}} \text{ pointed} \Rightarrow \mathcal{K}^* \text{ full-dimensional} \quad (315)$$

Given closed convex cone  $\mathcal{K} \subset \mathbb{R}^n$ , [58, §3.3 exer.20]<sup>2.65</sup>

$$\begin{array}{ccc} \mathcal{K} \text{ is pointed} & \Leftrightarrow & \text{full-dimensional} \\ \text{full-dimensional} & \Leftrightarrow & \text{pointed} \end{array} \quad \text{is } \mathcal{K}^* \quad (316)$$

The dual, of a closed convex cone not full-dimensional, contains a line. (§2.7.2.1.2)

- (vector sum) [343, thm.3.8] For convex cones  $\mathcal{K}_1$  and  $\mathcal{K}_2$

$$\mathcal{K}_1 + \mathcal{K}_2 = \text{conv}(\mathcal{K}_1 \cup \mathcal{K}_2) \quad (317)$$

is a convex cone.

- (dual vector-sum) [343, §16.4.2] [122, §4.6] For convex cones  $\mathcal{K}_1$  and  $\mathcal{K}_2$

$$\mathcal{K}_1^* \cap \mathcal{K}_2^* = (\mathcal{K}_1 + \mathcal{K}_2)^* = (\mathcal{K}_1 \cup \mathcal{K}_2)^* \quad (318)$$

- (closure of vector sum of duals)<sup>2.66</sup> For closed convex cones  $\mathcal{K}_1$  and  $\mathcal{K}_2$

$$(\mathcal{K}_1 \cap \mathcal{K}_2)^* = \overline{\mathcal{K}_1^* + \mathcal{K}_2^*} = \overline{\text{conv}(\mathcal{K}_1^* \cup \mathcal{K}_2^*)} \quad (319)$$

[366, p.96] where operation closure becomes superfluous under sufficient condition  $\mathcal{K}_1 \cap \text{intr } \mathcal{K}_2 \neq \emptyset$  [58, §3.3 exer.16, §4.1 exer.7].

- (Krein-Rutman) Given closed convex cones  $\mathcal{K}_1 \subseteq \mathbb{R}^m$  and  $\mathcal{K}_2 \subseteq \mathbb{R}^n$  and any linear map  $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , then provided  $\text{intr } \mathcal{K}_1 \cap A\mathcal{K}_2 \neq \emptyset$  [58, §3.3.13, confer §4.1 exer.9]

$$(A^{-1}\mathcal{K}_1 \cap \mathcal{K}_2)^* = A^T \mathcal{K}_1^* + \mathcal{K}_2^* \quad (320)$$

where dual of cone  $\mathcal{K}_1$  is with respect to its ambient space  $\mathbb{R}^m$  and dual of cone  $\mathcal{K}_2$  is with respect to  $\mathbb{R}^n$ , where  $A^{-1}\mathcal{K}_1$  denotes inverse image (§2.1.9.0.1) of  $\mathcal{K}_1$  under mapping  $A$ , and where  $A^T$  denotes adjoint operator. The particularly important case  $\mathcal{K}_2 = \mathbb{R}^n$  is easy to show: for  $A^T A = I$

$$\begin{aligned} (A^T \mathcal{K}_1)^* &\triangleq \{y \in \mathbb{R}^n \mid x^T y \geq 0 \ \forall x \in A^T \mathcal{K}_1\} \\ &= \{y \in \mathbb{R}^n \mid (A^T z)^T y \geq 0 \ \forall z \in \mathcal{K}_1\} \\ &= \{A^T w \mid z^T w \geq 0 \ \forall z \in \mathcal{K}_1\} \\ &= A^T \mathcal{K}_1^* \end{aligned} \quad (321)$$

- $\mathcal{K}$  is proper if and only if  $\mathcal{K}^*$  is proper.
- $\mathcal{K}$  is polyhedral if and only if  $\mathcal{K}^*$  is polyhedral. [366, §2.8]
- $\mathcal{K}$  is simplicial if and only if  $\mathcal{K}^*$  is simplicial. (§2.13.10.2) A simplicial cone and its dual are proper polyhedral cones (Figure 68, Figure 58), but not the converse.
- $\mathcal{K} \boxplus -\mathcal{K}^* = \mathbb{R}^n \Leftrightarrow \mathcal{K}$  is closed and convex. (2224)
- Any direction in a proper cone  $\mathcal{K}$  is normal to a hyperplane separating  $\mathcal{K}$  from  $-\mathcal{K}^*$ .

<sup>2.65</sup>  $\mathcal{K}^*$  is full-dimensional iff  $\mathcal{K}^* - \mathcal{K}^* = \mathbb{R}^n$ .

<sup>2.66</sup> These parallel analogous results for subspaces  $\mathcal{R}_1, \mathcal{R}_2 \subseteq \mathbb{R}^n$ ; [122, §4.6]

$$\begin{aligned} (\mathcal{R}_1 + \mathcal{R}_2)^\perp &= \mathcal{R}_1^\perp \cap \mathcal{R}_2^\perp \\ (\mathcal{R}_1 \cap \mathcal{R}_2)^\perp &= \overline{\mathcal{R}_1^\perp + \mathcal{R}_2^\perp} \end{aligned}$$

$\mathcal{R}^{\perp\perp} = \mathcal{R}$  for any subspace  $\mathcal{R}$ .

### 2.13.2 Abstractions of *Farkas' lemma*

**2.13.2.0.1 Corollary.** *Generalized inequality and membership relation.* [225, §A.4.2]  
Let  $\mathcal{K}$  be any closed convex cone and  $\mathcal{K}^*$  its dual, and let  $x$  and  $y$  belong to a vector space  $\mathbb{R}^n$ . Then

$$y \in \mathcal{K}^* \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } x \in \mathcal{K} \quad (322)$$

which is, merely, a statement of fact by definition of dual cone (300). By closure we have conjugation: [343, thm.14.1]

$$x \in \mathcal{K} \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } y \in \mathcal{K}^* \quad (323)$$

which may be regarded as a simple translation of *Farkas' lemma* [151] as in [343, §22] to the language of convex cones, and a generalization of the well-known Cartesian cone fact

$$x \succeq 0 \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } y \succeq 0 \quad (324)$$

for which implicitly  $\mathcal{K} = \mathcal{K}^* = \mathbb{R}_+^n$  the nonnegative orthant.

*Membership relation* (323) is often written instead as *dual generalized inequalities*, when  $\mathcal{K}$  and  $\mathcal{K}^*$  are pointed closed convex cones,

$$\begin{matrix} x \succeq 0 \\ \mathcal{K} \end{matrix} \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } \begin{matrix} y \succeq 0 \\ \mathcal{K}^* \end{matrix} \quad (325)$$

meaning, coordinates for biorthogonal expansion of  $x$  (§2.13.8.1.2, §2.13.9) [406] must be nonnegative when  $x$  belongs to  $\mathcal{K}$ . Conjugating,

$$\begin{matrix} y \succeq 0 \\ \mathcal{K}^* \end{matrix} \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } \begin{matrix} x \succeq 0 \\ \mathcal{K} \end{matrix} \quad (326)$$

◊

When pointed closed convex cone  $\mathcal{K}$  is not polyhedral, coordinate axes for biorthogonal expansion asserted by the corollary are taken from extreme directions of  $\mathcal{K}$ ; expansion is assured by *Carathéodory's theorem* (§E.6.4.1.1).

We presume, throughout, the obvious:

$$\begin{aligned} x \in \mathcal{K} &\Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } y \in \mathcal{K}^* & (323) \\ &\Leftrightarrow \\ x \in \mathcal{K} &\Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } y \in \mathcal{K}^*, \|y\|=1 & (327) \end{aligned}$$

#### 2.13.2.0.2 Exercise. Dual generalized inequalities.

Test Corollary 2.13.2.0.1 and (327) graphically on the two-dimensional polyhedral cone and its dual in Figure 60. ▼

(confer §2.7.2.2) When pointed closed convex cone  $\mathcal{K}$  is implicit from context:

$$\begin{aligned} x \succeq 0 &\Leftrightarrow x \in \mathcal{K} \\ x \succ 0 &\Leftrightarrow x \in \text{rel intr } \mathcal{K} \end{aligned} \quad (328)$$

Strict inequality  $x \succ 0$  means coordinates for biorthogonal expansion of  $x$  must be positive when  $x$  belongs to  $\text{rel intr } \mathcal{K}$ . Strict membership relations are useful; e.g., for any proper cone<sup>2.67</sup>  $\mathcal{K}$  and its dual  $\mathcal{K}^*$

$$x \in \text{intr } \mathcal{K} \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } y \in \mathcal{K}^*, y \neq \mathbf{0} \quad (329)$$

$$x \in \mathcal{K}, x \neq \mathbf{0} \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } y \in \text{intr } \mathcal{K}^* \quad (330)$$

<sup>2.67</sup>An open cone  $\mathcal{K}$  is admitted to (329) and (332) by (19).

Conjugating, we get the dual relations:

$$y \in \text{intr } \mathcal{K}^* \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } x \in \mathcal{K}, x \neq \mathbf{0} \quad (331)$$

$$y \in \mathcal{K}^*, y \neq \mathbf{0} \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } x \in \text{intr } \mathcal{K} \quad (332)$$

Boundary-membership relations for proper cones are also useful:

$$x \in \partial \mathcal{K} \Leftrightarrow \exists y \neq \mathbf{0} \ni \langle y, x \rangle = 0, y \in \mathcal{K}^*, x \in \mathcal{K} \quad (333)$$

$$y \in \partial \mathcal{K}^* \Leftrightarrow \exists x \neq \mathbf{0} \ni \langle y, x \rangle = 0, x \in \mathcal{K}, y \in \mathcal{K}^* \quad (334)$$

which are consistent; *e.g.*  $x \in \partial \mathcal{K} \Leftrightarrow x \notin \text{intr } \mathcal{K}$  **and**  $x \in \mathcal{K}$ .

#### 2.13.2.0.3 Example. Linear inequality. [374, §4] (confer §2.13.6.1.1)

Consider a given matrix  $A$  and closed convex cone  $\mathcal{K}$ . By membership relation we have

$$\begin{aligned} Ay \in \mathcal{K}^* &\Leftrightarrow x^T A y \geq 0 \quad \forall x \in \mathcal{K} \\ &\Leftrightarrow y^T z \geq 0 \quad \forall z \in \{A^T x \mid x \in \mathcal{K}\} \\ &\Leftrightarrow y \in \{A^T x \mid x \in \mathcal{K}\}^* \end{aligned} \quad (335)$$

This implies

$$\{y \mid Ay \in \mathcal{K}^*\} = \{A^T x \mid x \in \mathcal{K}\}^* \quad (336)$$

When  $\mathcal{K}$  is the *selfdual* nonnegative orthant (§2.13.6.1), for example, then

$$\{y \mid Ay \succeq 0\} = \{A^T x \mid x \succeq 0\}^* \quad (337)$$

and the dual relation

$$\{y \mid Ay \succeq 0\}^* = \{A^T x \mid x \succeq 0\} \quad (338)$$

comes by a theorem of Weyl (p.117) that yields closedness for any vertex-description of a polyhedral cone.  $\square$

#### 2.13.2.1 Null certificate, Theorem of the alternative

If in particular  $x_p \notin \mathcal{K}$  a closed convex cone, then construction in Figure 59b suggests there exists a supporting hyperplane (having inward-normal belonging to dual cone  $\mathcal{K}^*$ ) separating  $x_p$  from  $\mathcal{K}$ ; indeed, (323)

$$x_p \notin \mathcal{K} \Leftrightarrow \exists y \in \mathcal{K}^* \ni \langle y, x_p \rangle < 0 \quad (339)$$

Existence of any one such  $y$  is a certificate of null membership. From a different perspective,

$$\begin{aligned} x_p &\in \mathcal{K} \\ \text{or in the alternative} \\ \exists y \in \mathcal{K}^* \ni \langle y, x_p \rangle &< 0 \end{aligned} \quad (340)$$

By *alternative* is meant: these two systems are incompatible; one system is feasible while the other is not.

### 2.13.2.1.1 Example. Theorem of the alternative for linear inequality.

Myriad alternative systems of linear inequality can be explained in terms of pointed closed convex cones and their duals.

Beginning from the simplest Cartesian dual generalized inequalities (324) (with respect to nonnegative orthant  $\mathbb{R}_+^m$ ),

$$y \succeq 0 \Leftrightarrow x^T y \geq 0 \text{ for all } x \succeq 0 \quad (341)$$

Given  $A \in \mathbb{R}^{n \times m}$ , we make vector substitution  $y \leftarrow A^T y$

$$A^T y \succeq 0 \Leftrightarrow x^T A^T y \geq 0 \text{ for all } x \succeq 0 \quad (342)$$

Introducing a new vector by calculating  $b \triangleq Ax$  we get

$$A^T y \succeq 0 \Leftrightarrow b^T y \geq 0, \quad b = Ax \text{ for all } x \succeq 0 \quad (343)$$

By complementing sense of the scalar inequality:

$$\begin{aligned} A^T y \succeq 0 \\ \text{or in the alternative} \\ b^T y < 0, \quad \exists b = Ax, \quad x \succeq 0 \end{aligned} \quad (344)$$

If one system has a solution, then the other does not; define a convex cone  $\mathcal{K} = \{y \mid A^T y \succeq 0\}$ , then  $y \in \mathcal{K}$  or in the alternative  $y \notin \mathcal{K}$ .

Scalar inequality  $b^T y < 0$  is movable to the other side of alternative (344), but that requires some explanation: From results in Example 2.13.2.0.3, the dual cone is  $\mathcal{K}^* = \{Ax \mid x \succeq 0\}$ . From (323) we have

$$\begin{aligned} y \in \mathcal{K} &\Leftrightarrow b^T y \geq 0 \text{ for all } b \in \mathcal{K}^* \\ A^T y \succeq 0 &\Leftrightarrow b^T y \geq 0 \text{ for all } b \in \{Ax \mid x \succeq 0\} \end{aligned} \quad (345)$$

Given some  $b$  vector and  $y \in \mathcal{K}$ , then  $b^T y < 0$  can only mean  $b \notin \mathcal{K}^*$ . An alternative system is therefore simply  $b \in \mathcal{K}^*$ : [225, p.59] (Farkas/Tucker)

$$\begin{aligned} A^T y \succeq 0, \quad b^T y < 0 \\ \text{or in the alternative} \\ b = Ax, \quad x \succeq 0 \end{aligned} \quad (346)$$

Geometrically this means: either vector  $b$  belongs to convex cone  $\mathcal{K}^*$  or it does not. When  $b \notin \mathcal{K}^*$ , then there is a vector  $y \in \mathcal{K}$  normal to a hyperplane separating point  $b$  from cone  $\mathcal{K}^*$ .

For another example, from membership relation (322) with affine transformation of dual variable we may write: Given  $A \in \mathbb{R}^{n \times m}$  and  $b \in \mathbb{R}^n$

$$b - Ay \in \mathcal{K}^* \Leftrightarrow x^T(b - Ay) \geq 0 \quad \forall x \in \mathcal{K} \quad (347)$$

$$A^T x = \mathbf{0}, \quad b - Ay \in \mathcal{K}^* \Rightarrow x^T b \geq 0 \quad \forall x \in \mathcal{K} \quad (348)$$

From membership relation (347), conversely, suppose we allow any  $y \in \mathbb{R}^m$ . Then because  $-x^T A y$  is unbounded below,  $x^T(b - Ay) \geq 0$  implies  $A^T x = \mathbf{0}$ : for  $y \in \mathbb{R}^m$

$$A^T x = \mathbf{0}, \quad b - Ay \in \mathcal{K}^* \Leftrightarrow x^T(b - Ay) \geq 0 \quad \forall x \in \mathcal{K} \quad (349)$$

*In toto,*

$$b - Ay \in \mathcal{K}^* \Leftrightarrow x^T b \geq 0, \quad A^T x = \mathbf{0} \quad \forall x \in \mathcal{K} \quad (350)$$

Vector  $x$  belongs to cone  $\mathcal{K}$  but is also constrained to lie in a subspace of  $\mathbb{R}^n$  specified by an intersection of hyperplanes through the origin  $\{x \in \mathbb{R}^n \mid A^T x = \mathbf{0}\}$ . From this, alternative systems of generalized inequality with respect to pointed closed convex cones  $\mathcal{K}$  and  $\mathcal{K}^*$

$$\begin{array}{c} Ay \preceq b \\ \mathcal{K}^* \\ \text{or in the alternative} \\ x^T b < 0, \quad A^T x = \mathbf{0}, \quad x \succeq_0 0 \end{array} \quad (351)$$

derived from (350) simply by taking the complementary sense of the inequality in  $x^T b$ . These two systems are alternatives; if one system has a solution, then the other does not.<sup>2.68</sup> [343, p.201]

By invoking a strict membership relation between proper cones (329), we can construct a more exotic alternative strengthened by demand for an interior point;

$$b - Ay \succ 0 \Leftrightarrow x^T b > 0, \quad A^T x = \mathbf{0} \quad \forall x \succeq_0 0, \quad x \neq \mathbf{0} \quad (352)$$

From this, alternative systems of generalized inequality [65, pp.50,54,262]

$$\begin{array}{c} Ay \prec b \\ \mathcal{K}^* \\ \text{or in the alternative} \\ x^T b \leq 0, \quad A^T x = \mathbf{0}, \quad x \succeq_0 0, \quad x \neq \mathbf{0} \end{array} \quad (353)$$

derived from (352) by taking complementary sense of the inequality in  $x^T b$ .

And from this, alternative systems with respect to the nonnegative orthant attributed to Gordan in 1873: [183] [58, §2.2] substituting  $A \leftarrow A^T$  and setting  $b = \mathbf{0}$

$$\begin{array}{c} A^T y \prec 0 \\ \text{or in the alternative} \\ Ax = \mathbf{0}, \quad x \succeq 0, \quad \|x\|_1 = 1 \end{array} \quad (354)$$

Ben-Israel collects related results from Farkas, Motzkin, Gordan, and Stiemke in *Motzkin transposition theorem*. [34]  $\square$

### 2.13.3 Optimality condition

(confer §2.13.11.1) The general first-order necessary and sufficient condition for optimality of solution  $x^*$  to a minimization problem ((306p) for example) with real differentiable convex objective function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  is [342, §3]

$$\nabla f(x^*)^T (x - x^*) \geq 0 \quad \forall x \in \mathcal{C}, \quad x^* \in \mathcal{C} \quad (355)$$

---

<sup>2.68</sup>If solutions at  $\pm\infty$  are disallowed, then the alternative systems become instead *mutually exclusive* with respect to nonpolyhedral cones. Simultaneous infeasibility of the two systems is not precluded by mutual exclusivity; called a *weak alternative*. Ye provides an example illustrating simultaneous infeasibility with respect to the positive semidefinite cone:  $x \in \mathbb{S}^2$ ,  $y \in \mathbb{R}$ ,  $A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ , and  $b = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$  where  $x^T b$  means  $\langle x, b \rangle$ . A better strategy than disallowing solutions at  $\pm\infty$  is to demand an interior point as in (353) or Lemma 4.2.1.1.2. Then question of simultaneous infeasibility is moot.

where  $\mathcal{C}$  is a convex *feasible set*,<sup>2.69</sup> and where  $\nabla f(x^*)$  is the *gradient* (§3.6) of  $f$  with respect to  $x$  evaluated at  $x^*$ . In words, negative gradient is normal to a hyperplane supporting the feasible set at a point of optimality. (Figure 71)

Direct solution to *variation inequality* (355), instead of the corresponding minimization, spawned from *calculus of variations*. [280, p.178] [150, p.37] One solution method solves an equivalent fixed point-of-projection problem

$$x = P_{\mathcal{C}}(x - \nabla f(x)) \quad (356)$$

that follows from a necessary and sufficient condition for projection on convex set  $\mathcal{C}$  (Theorem E.9.1.0.2)

$$P(x^* - \nabla f(x^*)) \in \mathcal{C}, \quad \langle x^* - \nabla f(x^*) - x^*, x - x^* \rangle \leq 0 \quad \forall x \in \mathcal{C} \quad (2208)$$

Proof of equivalence [410, *Complementarity problem*] is provided by Németh. Given minimum-distance projection problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && \frac{1}{2} \|x - y\|^2 \\ & \text{subject to} && x \in \mathcal{C} \end{aligned} \quad (357)$$

on convex feasible set  $\mathcal{C}$  for example, the equivalent fixed point problem

$$x = P_{\mathcal{C}}(x - \nabla f(x)) = P_{\mathcal{C}}(y) \quad (358)$$

is solved in one step.

In the unconstrained case ( $\mathcal{C} = \mathbb{R}^n$ ), optimality condition (355) reduces to the classical rule (p.194)

$$\nabla f(x^*) = \mathbf{0}, \quad x^* \in \text{dom } f \quad (359)$$

which can be inferred from the following application:

#### 2.13.3.0.1 Example. Optimality for equality-constrained problem.

Given a real differentiable convex function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  defined on domain  $\mathbb{R}^n$ , a wide full-rank matrix  $C \in \mathbb{R}^{p \times n}$ , and vector  $d \in \mathbb{R}^p$ , the convex optimization problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{subject to} && Cx = d \end{aligned} \quad (360)$$

is characterized by the well-known necessary and sufficient optimality condition [65, §4.2.3]

$$\nabla f(x^*) + C^T \nu = \mathbf{0} \quad (361)$$

where  $\nu \in \mathbb{R}^p$  is the eminent *Lagrange multiplier*. [341] [280, p.188] [257] In other words, solution  $x^*$  is optimal if and only if  $\nabla f(x^*)$  belongs to  $\mathcal{R}(C^T)$ .

Via membership relation, we now derive condition (361) from the general first-order condition for optimality (355): For problem (360)

$$\mathcal{C} \triangleq \{x \in \mathbb{R}^n \mid Cx = d\} = \{Z\xi + x_p \mid \xi \in \mathbb{R}^{n-\text{rank } C}\} \quad (362)$$

is the feasible set where  $Z \in \mathbb{R}^{n \times n-\text{rank } C}$  holds basis  $\mathcal{N}(C)$  columnar, and  $x_p$  is any particular solution to  $Cx = d$ . Since  $x^* \in \mathcal{C}$ , we arbitrarily choose  $x_p = x^*$  which yields an equivalent optimality condition

$$\nabla f(x^*)^T Z\xi \geq 0 \quad \forall \xi \in \mathbb{R}^{n-\text{rank } C} \quad (363)$$

---

<sup>2.69</sup> presumably nonempty set of all variable values satisfying all given problem constraints; the set of *feasible solutions*.

when substituted into (355). But this is simply half of a membership relation where the cone dual to  $\mathbb{R}^{n-\text{rank } C}$  is the origin in  $\mathbb{R}^{n-\text{rank } C}$ . We must therefore have

$$Z^T \nabla f(x^*) = \mathbf{0} \Leftrightarrow \nabla f(x^*)^T Z \xi \geq 0 \quad \forall \xi \in \mathbb{R}^{n-\text{rank } C} \quad (364)$$

meaning,  $\nabla f(x^*)$  must be orthogonal to  $\mathcal{N}(C)$ . These conditions

$$Z^T \nabla f(x^*) = \mathbf{0}, \quad Cx^* = d \quad (365)$$

are necessary and sufficient for optimality.  $\square$

### 2.13.4 Discretization of membership relation

#### 2.13.4.1 Dual halfspace-description

Halfspace-description of dual cone  $\mathcal{K}^*$  is equally simple as vertex-description

$$\mathcal{K} = \text{cone}(X) = \{Xa \mid a \succeq 0\} \subseteq \mathbb{R}^n \quad (105)$$

for corresponding closed convex cone  $\mathcal{K}$ : By definition (300), for  $X \in \mathbb{R}^{n \times N}$  as in (282), (confer (289))

$$\begin{aligned} \mathcal{K}^* &= \left\{ y \in \mathbb{R}^n \mid z^T y \geq 0 \text{ for all } z \in \mathcal{K} \right\} \\ &= \left\{ y \in \mathbb{R}^n \mid z^T y \geq 0 \text{ for all } z = Xa, a \succeq 0 \right\} \\ &= \left\{ y \in \mathbb{R}^n \mid a^T X^T y \geq 0 \text{ for all } a \succeq 0 \right\} \\ &= \left\{ y \in \mathbb{R}^n \mid X^T y \succeq 0 \right\} \end{aligned} \quad (366)$$

that follows from the *generalized inequality and membership corollary* (324). The semi-infinity of tests specified by all  $z \in \mathcal{K}$  has been reduced to a set of generators for  $\mathcal{K}$  constituting the columns of  $X$ ; *id est*, the test has been discretized.

Whenever cone  $\mathcal{K}$  is known to be closed and convex, the conjugate statement must also hold; *id est*, given any set of generators for dual cone  $\mathcal{K}^*$  arranged columnar in  $Y$ , then the consequent vertex-description of dual cone connotes a halfspace-description for  $\mathcal{K}$ : [366, §2.8]

$$\mathcal{K}^* = \{Ya \mid a \succeq 0\} \Leftrightarrow \mathcal{K}^{**} = \mathcal{K} = \{z \mid Y^T z \succeq 0\} \quad (367)$$

#### 2.13.4.2 First dual-cone formula

From these two results (366) and (367) we deduce a general principle:

- From any [*sic*] given vertex-description (105) of closed convex cone  $\mathcal{K}$ , halfspace-description (366) of the dual cone  $\mathcal{K}^*$  is immediate by matrix transposition. Conversely, from any given halfspace-description (289) of  $\mathcal{K}$ , dual vertex-description (367) is immediate. [343, p.122]

Various other converses are just a little trickier. (§2.13.10, §2.13.12)

We further deduce: For any polyhedral cone  $\mathcal{K}$ , the dual cone  $\mathcal{K}^*$  is also polyhedral and  $\mathcal{K}^{**} = \mathcal{K}$ . [366, p.56] For any pointed polyhedral cone  $\mathcal{K}$ , dual cone  $\mathcal{K}^*$  is full-dimensional. (315) (§2.13.7)

The *generalized inequality and membership corollary* is discretized in the following theorem inspired by (366) and (367):

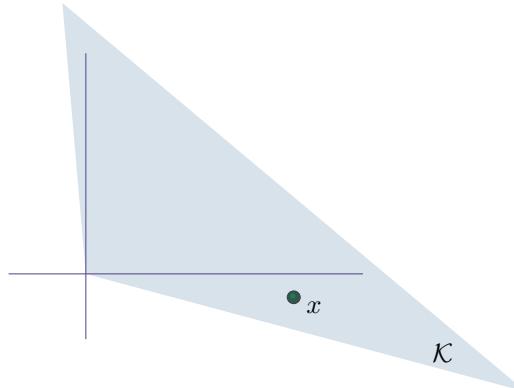


Figure 65:  $x \succeq 0$  with respect to  $\mathcal{K}$  but not with respect to nonnegative orthant  $\mathbb{R}_+^2$  (pointed convex cone  $\mathcal{K}$  drawn truncated).

**2.13.4.2.1 Theorem.** *Discretized membership.* (confer §2.13.2.0.1)<sup>2.70</sup>

Given any set of generators (§2.8.1.2) denoted by  $\mathcal{G}(\mathcal{K})$ , for closed convex cone  $\mathcal{K} \subseteq \mathbb{R}^n$ , and any set of generators denoted  $\mathcal{G}(\mathcal{K}^*)$  for its dual such that

$$\mathcal{K} = \text{cone } \mathcal{G}(\mathcal{K}), \quad \mathcal{K}^* = \text{cone } \mathcal{G}(\mathcal{K}^*) \quad (368)$$

then discretization of the *generalized inequality and membership corollary* (p.131) is necessary and sufficient for certifying cone membership: for  $x$  and  $y$  in vector space  $\mathbb{R}^n$

$$x \in \mathcal{K} \Leftrightarrow \langle \gamma^*, x \rangle \geq 0 \text{ for all } \gamma^* \in \mathcal{G}(\mathcal{K}^*) \quad (369)$$

$$y \in \mathcal{K}^* \Leftrightarrow \langle \gamma, y \rangle \geq 0 \text{ for all } \gamma \in \mathcal{G}(\mathcal{K}) \quad (370)$$

◇

**Proof.**  $\mathcal{K}^* = \{\mathcal{G}(\mathcal{K}^*)a \mid a \succeq 0\}$ .  $y \in \mathcal{K}^* \Leftrightarrow y = \mathcal{G}(\mathcal{K}^*)a \text{ for some } a \succeq 0$ .  $x \in \mathcal{K} \Leftrightarrow \langle y, x \rangle \geq 0 \forall y \in \mathcal{K}^* \Leftrightarrow \langle \mathcal{G}(\mathcal{K}^*)a, x \rangle \geq 0 \forall a \succeq 0$  (323).  $a \triangleq \sum_i \alpha_i e_i$  where  $e_i$  is the  $i^{\text{th}}$  member of a standard basis of possibly infinite cardinality.  $\langle \mathcal{G}(\mathcal{K}^*)a, x \rangle \geq 0 \forall a \succeq 0 \Leftrightarrow \sum_i \alpha_i \langle \mathcal{G}(\mathcal{K}^*)e_i, x \rangle \geq 0 \forall \alpha_i \geq 0 \Leftrightarrow \langle \mathcal{G}(\mathcal{K}^*)e_i, x \rangle \geq 0 \forall i$ . Conjugate relation (370) is similarly derived. ♦

**2.13.4.2.2 Exercise.** *Discretized dual generalized inequalities.*

Test Theorem 2.13.4.2.1 on Figure 60a using extreme directions there as generators.

▼

From the *discretized membership theorem* we may further deduce a more surgical description of closed convex cone that prescribes intersection of only a finite number of halfspaces for its construction when polyhedral: (Figure 59a)

$$\mathcal{K} = \{x \in \mathbb{R}^n \mid \langle \gamma^*, x \rangle \geq 0 \text{ for all } \gamma^* \in \mathcal{G}(\mathcal{K}^*)\} \quad (371)$$

$$\mathcal{K}^* = \{y \in \mathbb{R}^n \mid \langle \gamma, y \rangle \geq 0 \text{ for all } \gamma \in \mathcal{G}(\mathcal{K})\} \quad (372)$$

<sup>2.70</sup>Stated in [22, §1] without proof for pointed closed convex case.

#### 2.13.4.2.3 Exercise. Partial order induced by orthant.

When comparison is with respect to the nonnegative orthant  $\mathcal{K} = \mathbb{R}_+^n$ , then from the *discretized membership theorem* it directly follows:

$$x \preceq z \Leftrightarrow x_i \leq z_i \quad \forall i \quad (373)$$

Generate simple counterexamples demonstrating that this equivalence with entrywise inequality holds only when the underlying cone inducing partial order is the nonnegative orthant; *e.g.*, explain Figure 65.  $\blacktriangledown$

#### 2.13.4.2.4 Example. Boundary membership to polyhedral cone.

For a polyhedral cone, test (333) of boundary membership can be formulated as a linear program. Say proper polyhedral cone  $\mathcal{K}$  is specified completely by generators that are arranged columnar in

$$X = [\Gamma_1 \ \cdots \ \Gamma_N] \in \mathbb{R}^{n \times N} \quad (282)$$

*id est*,  $\mathcal{K} = \{Xa \mid a \succeq 0\}$  (105). Then boundary-membership relation for proper cone

$$c \in \partial\mathcal{K} \Leftrightarrow \exists y \neq \mathbf{0} \ni \langle y, c \rangle = 0, \quad y \in \mathcal{K}^*, \quad c \in \mathcal{K} \quad (333)$$

may be expressed<sup>2.71</sup>

$$\begin{aligned} & \underset{a, y}{\text{find}} \quad y \neq \mathbf{0} \\ & \text{subject to} \quad c^T y = 0 \\ & \quad X^T y \succeq 0 \\ & \quad Xa = c \\ & \quad a \succeq 0 \end{aligned} \quad (374)$$

This linear feasibility problem has a solution iff  $c \in \partial\mathcal{K}$ . If membership  $c \in \mathcal{K}$  is known *a priori*, then variable  $a$  becomes redundant. This method assumes a full-dimensional cone.

We may adapt (374) to cones, contained wholly in a subspace, by introducing affine hull  $\text{aff } \mathcal{K}$  into the program; *id est*, given  $c \in \mathcal{K}$

$$\begin{aligned} & \underset{h, y}{\text{find}} \quad y \neq \mathbf{0} \\ & \text{subject to} \quad c^T y = 0 \\ & \quad X^T y \succeq 0 \\ & \quad y = [\mathbf{0} \ X]h \\ & \quad h^T \mathbf{1} = 1 \end{aligned} \quad (375)$$

Now this linear feasibility problem has solution iff  $c \in \text{rel } \partial\mathcal{K}$ . This adaptation is necessary to determine boundary membership to a cone that is not full-dimensional.  $\square$

### 2.13.5 Smallest face of closed convex cone

Given nonempty convex subset  $\mathcal{C}$  of a convex set  $\mathcal{K}$ , the *smallest face* of  $\mathcal{K}$  containing  $\mathcal{C}$  is equivalent to intersection of all faces of  $\mathcal{K}$  that contain  $\mathcal{C}$ . [343, p.164] By (313), membership relation (333) means that each and every point on boundary  $\partial\mathcal{K}$  of proper cone  $\mathcal{K}$  belongs to a hyperplane supporting  $\mathcal{K}$  whose normal  $y$  belongs to dual cone  $\mathcal{K}^*$ . It follows that the smallest face  $\mathcal{F}$ , containing  $\mathcal{C} \subset \partial\mathcal{K} \subset \mathbb{R}^n$  on boundary of proper cone  $\mathcal{K}$ , is the intersection of all hyperplanes containing  $\mathcal{C}$  whose normals are in  $\mathcal{K}^*$ ;

$$\mathcal{F}(\mathcal{K} \supset \mathcal{C}) = \{x \in \mathcal{K} \mid x \perp \mathcal{K}^* \cap \mathcal{C}^\perp\} \quad (376)$$

where

$$\mathcal{C}^\perp \triangleq \{y \in \mathbb{R}^n \mid \langle z, y \rangle = 0 \quad \forall z \in \mathcal{C}\} \quad (377)$$

When  $\mathcal{C} \cap \text{intr } \mathcal{K} \neq \emptyset$  then  $\mathcal{F}(\mathcal{K} \supset \mathcal{C}) = \mathcal{K}$ .

---

<sup>2.71</sup>A clumsy but sure convex method, for determining whether nonzero  $y \in \mathbb{R}^n$  exists, is to constrain each entry  $y_i = \pm 1$ ,  $i = 1 \dots n$  individually until one of them (if any) becomes feasible.

**2.13.5.0.1 Example.** *Finding smallest face of cone.*

Suppose polyhedral cone  $\mathcal{K}$  is completely specified by generators arranged columnar in

$$X = [\Gamma_1 \ \cdots \ \Gamma_N] \in \mathbb{R}^{n \times N} \quad (282)$$

*id est,*

$$\mathcal{K} = \text{cone } X = \{Xa \mid a \succeq 0\} \subseteq \mathbb{R}^n \quad (105)$$

To find its smallest face  $\mathcal{F}(\mathcal{K} \ni c)$  containing a given point  $c \in \mathcal{K}$ , by the *discretized membership theorem* 2.13.4.2.1, it is necessary and sufficient to find generators for the smallest face. We may do so one generator at a time:

Consider generator  $\Gamma_i$ . If there exists a vector  $z \in \mathcal{K}^*$  in the dual cone that is orthogonal to  $c$  but not to  $\Gamma_i$ , then  $\Gamma_i$  cannot belong to the smallest face of  $\mathcal{K}$  containing  $c$ . Such a vector  $z$  can be realized by a linear feasibility problem: given  $c \in \mathcal{K}$

$$\begin{aligned} &\text{find} && z \in \mathbb{R}^n \\ &\text{subject to} && c^T z = 0 \\ &&& X^T z \succeq 0 \\ &&& \Gamma_i^T z = 1, \quad i \in \{1 \dots N\} \end{aligned} \quad (378)$$

If there exists a solution  $z$  for which  $\Gamma_i^T z = 1$ , then

$$\Gamma_i \not\subset \mathcal{K}^* \cap c^\perp = \{z \in \mathbb{R}^n \mid X^T z \succeq 0, c^T z = 0\} \quad (379)$$

so  $\Gamma_i \notin \mathcal{F}(\mathcal{K} \ni c)$ ; solution  $z$  is a certificate of null membership. If problem (378) is infeasible for generator  $\Gamma_i \in \mathcal{K}$ , conversely, then  $\Gamma_i \in \mathcal{F}(\mathcal{K} \ni c)$  by (376) and (366) because  $\Gamma_i \perp \mathcal{K}^* \cap c^\perp$ ; in that case,  $\Gamma_i$  is a generator of  $\mathcal{F}(\mathcal{K} \ni c)$ .

Since the constant in constraint  $\Gamma_i^T z = 1$  is arbitrary positively, then there is correspondence between (378) and (351) [sic] admitting an alternative to linear feasibility problem (378): for a given point  $c \in \mathcal{K}$

$$\begin{aligned} &\text{find} && a, \mu \\ &\text{subject to} && \mu c - \Gamma_i = Xa, \quad i \in \{1 \dots N\} \\ &&& a \succeq 0 \end{aligned} \quad (380)$$

Now if this problem is feasible (bounded) for generator  $\Gamma_i \in \mathcal{K}$ , then (378) is infeasible and  $\Gamma_i \in \mathcal{F}(\mathcal{K} \ni c)$  is a generator of the smallest face of  $\mathcal{K}$  that contains  $c$ .

When finding a smallest face via (378) or (380), generators of  $\mathcal{K}$  in matrix  $X$  may not be diminished in number (by discarding columns) until all generators of the smallest face have been identified. Diminishing column space is a form of *presolving*; it is equivalent to a proof that coefficient  $a_i$  can only be 0 in

$$\begin{aligned} &\text{find} && a \in \mathbb{R}^N \\ &\text{subject to} && c = Xa \\ &&& a \succeq 0 \end{aligned} \quad (381)$$

because the corresponding  $\Gamma_i$  does not belong to the smallest face of  $\mathcal{K}$  that contains  $c$ .

*en masse*

Solving (378) or (380)  $N$  times can be computationally intensive if number of columns is large. If there were fewer than  $N$  unique<sup>2.72</sup> vectors  $z$  in the dual cone that satisfy (378), then it would be more economical to find them instead of testing columns of  $X$  one by

---

<sup>2.72</sup> unique in the same sense as for eigenvectors in Definition A.5.0.1.1.

one. We propose finding generators  $\{\Gamma_i\}$  *en masse* for the smallest face of polyhedral cone  $\mathcal{K}$ , containing  $c$ , by solving a variant of (378): assuming  $c \in \mathcal{K}$

$$\begin{aligned} & \underset{z \in \mathbb{R}^n}{\text{maximize}} \quad \mathbf{1}^T X^T z \\ & \text{subject to} \quad c^T z = 0 \\ & \quad \mathbf{1} \succeq X^T z \succeq 0 \\ & \quad -\mathbf{1} \preceq z \preceq \mathbf{1} \end{aligned} \tag{382}$$

where bounding to  $\mathbf{1}$  precludes an unbounded objective or variable and insures that maximization is democratic over all rows of  $X^T$ . For all  $\Gamma_i$  corresponding to an *inactive* inequality  $\{\Gamma_i^T z^* > 0\}$ , optimal solution  $z^*$  is a certificate of their null membership to  $\mathcal{F}(\mathcal{K} \ni c)$ . One way to find the set of all optimal  $\{z^*\}$  is to solve (382) recursively; *id est*, columns of  $X$ , corresponding to inactive inequalities, are deprecated (discarded) before solving (382) again. Recursion continues until the inequality constraint becomes completely active below:  $X^T z^* = \mathbf{0}$ . Surviving columns of  $X$  comprise a superset containing generators for  $\mathcal{F}(\mathcal{K} \ni c)$ .  $\square$

#### 2.13.5.0.2 Exercise. Optimality of null membership en masse.

Prove that solving (382) recursively is not equivalent to solving (378) for  $i=1 \dots N$ . While fast, (382) is suboptimal eliminating a proper though substantial subset (if not equal in number) of columns eliminated by (378).  $\blacktriangledown$

#### 2.13.5.0.3 Exercise. Finding smallest face for broader class of convex cone.

Show how algorithms (378) (380) (382) apply more broadly; *id est*, full-dimensionality<sup>2.73</sup> can be unnecessary.  $\blacktriangledown$

#### 2.13.5.0.4 Exercise. Finding smallest face by alternative system.

Derive (380) from (378).<sup>2.74</sup> What is a variant of (380) that finds generators  $\{\Gamma_i\}$  *en masse* for the smallest face of  $\mathcal{K}$  containing  $c$ .  $\blacktriangledown$

#### 2.13.5.0.5 Exercise. Smallest face of positive semidefinite cone.

Derive (225) from (376).  $\blacktriangledown$

#### 2.13.5.0.6 Exercise. Deprecation by column discard.

Explain why, when finding a smallest face via (378), generators of  $\mathcal{K}$  may be discarded only after all generators of the smallest face have been identified. Then explain why generators may be discarded *en masse* prior to finding a smallest face via (382).  $\blacktriangledown$

#### 2.13.5.0.7 Exercise. Elegantly ascertain membership to cone boundary.

Invent a convex objective or single convex constraint that insures finding a nonzero  $y$  if it exists in the feasible set of vectors  $y$  for problem (374).  $\blacktriangledown$

#### 2.13.5.0.8 Exercise. Constraints en masse.

Consider introducing a constraint  $\xi^T a = \phi$ , on coefficient vector  $a$  from (105), to *en masse* algorithm (382) for finding smallest face of polyhedral cone  $\mathcal{K}$  containing point  $c$ . Show that to be accomplished by merely augmenting the  $X$  matrix, as in  $[X^T \ \xi]$ , and the  $c$  vector as in  $[c^T \ \phi^T]$  given  $\xi$  and  $\phi$ .  $\blacktriangledown$

<sup>2.73</sup>Hint: A hyperplane, with normal in  $\mathcal{K}^*$ , containing cone  $\mathcal{K}$  is admissible; *e.g.*, Figure 44. More importantly, as in (375), full-dimensionality is obviated by application of affine hull of  $\mathcal{K}$ .

<sup>2.74</sup>Hint: Recall dual closed convex cone description (367), then swap cone  $\mathcal{K}$  with its dual in (351).

### 2.13.6 Dual PSD cone and generalized inequality

The *dual positive semidefinite cone*  $\mathcal{K}^*$  is confined to  $\mathbb{S}^M$  by convention:

$$\mathbb{S}_+^{M^*} \triangleq \{Y \in \mathbb{S}^M \mid \langle Y, X \rangle \geq 0 \text{ for all } X \in \mathbb{S}_+^M\} = \mathbb{S}_+^M \quad (383)$$

The positive semidefinite cone is selfdual in the ambient space of symmetric matrices [65, exmp.2.24] [40] [221, §II];  $\mathcal{K} = \mathcal{K}^*$ .

Dual generalized inequalities with respect to the positive semidefinite cone in the ambient space of symmetric matrices can therefore be simply stated: (Fejér)

$$X \succeq 0 \Leftrightarrow \text{tr}(Y^T X) \geq 0 \text{ for all } Y \succeq 0 \quad (384)$$

Membership to this cone can be determined in the isometrically isomorphic Euclidean space  $\mathbb{R}^{M^2}$  via (38). (§2.2.1) By the two interpretations in §2.13.1, positive semidefinite matrix  $Y$  can be interpreted as inward-normal to a hyperplane supporting the positive semidefinite cone.

The fundamental statement of positive semidefiniteness,  $y^T X y \geq 0 \forall y$  (§A.3.0.0.1), evokes a particular instance of these dual generalized inequalities (384):

$$X \succeq 0 \Leftrightarrow \langle yy^T, X \rangle \geq 0 \quad \forall yy^T(\succeq 0) \quad (1596)$$

Discretization (§2.13.4.2.1) allows replacement of positive semidefinite matrices  $Y$  with this minimal set of generators comprising the extreme directions of the positive semidefinite cone (§2.9.2.7).

#### 2.13.6.1 selfdual cones

From (133) (a consequence of the *halfspaces theorem*, §2.4.1.1.1), where the only finite value of the support function for a convex cone is 0 [225, §C.2.3.1], or from discretized definition (372) of the dual cone we get a rather self evident characterization of selfdual cones:

$$\mathcal{K} = \mathcal{K}^* \Leftrightarrow \mathcal{K} = \bigcap_{\gamma \in \mathcal{G}(\mathcal{K})} \{y \mid \gamma^T y \geq 0\} \quad (385)$$

In words: Cone  $\mathcal{K}$  is *selfdual* iff its own extreme directions are inward-normals to a (minimal) set of hyperplanes bounding halfspaces whose intersection constructs it. This means each extreme direction of  $\mathcal{K}$  is normal to a hyperplane exposing one of its own faces; a necessary but insufficient condition for selfdualness (Figure 66, for example).

Selfdual cones are necessarily full-dimensional. [31, §I] Their most prominent representatives are the orthants (Cartesian cones), the positive semidefinite cone  $\mathbb{S}_+^M$  in the ambient space of symmetric matrices (383), and Lorentz cone (181) [21, §II.A] [65, exmp.2.25]. In three dimensions, a plane containing the axis of revolution of a selfdual cone (and the origin) will produce a *slice* whose boundary makes a right angle.

##### 2.13.6.1.1 Example. Linear matrix inequality.

(confer §2.13.2.0.3)

Consider a peculiar vertex-description for a convex cone  $\mathcal{K}$  defined over a positive semidefinite cone (instead of a nonnegative orthant as in definition (105)): for  $X \in \mathbb{S}_+^n$

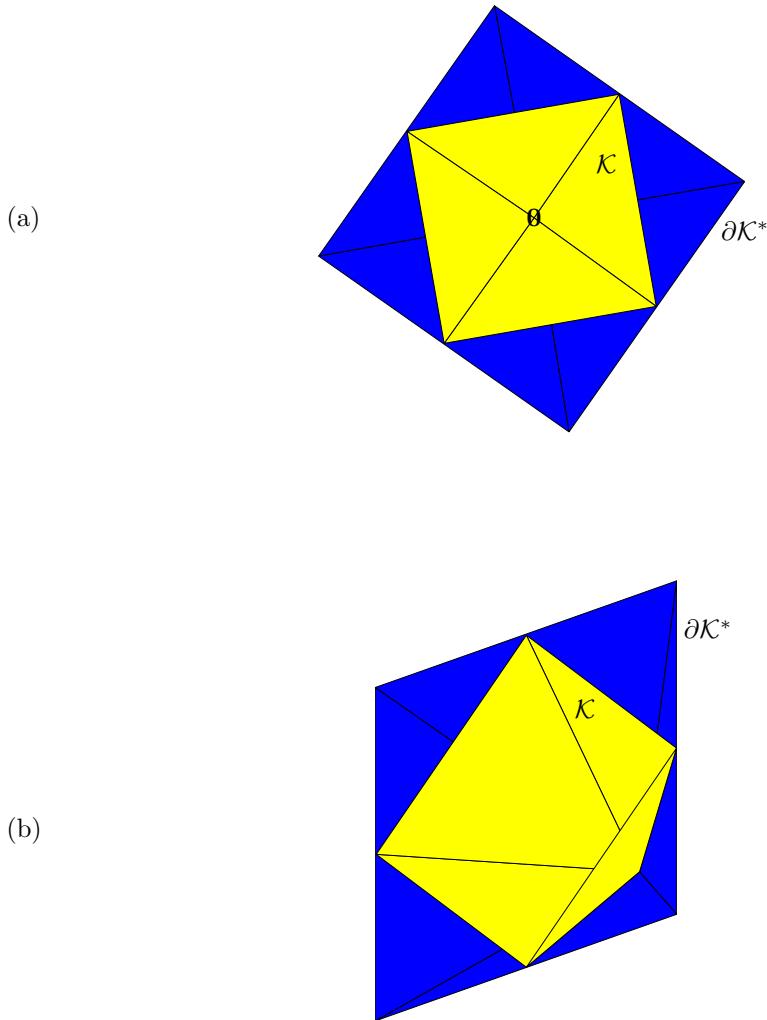


Figure 66: Two (truncated) views of a polyhedral cone  $\mathcal{K}$  and its dual in  $\mathbb{R}^3$ . Each of four extreme directions from  $\mathcal{K}$  belongs to a face of dual cone  $\mathcal{K}^*$ . Cone  $\mathcal{K}$ , shrouded by its dual, is symmetrical about its axis of revolution. Each pair of diametrically opposed extreme directions from  $\mathcal{K}$  makes a right angle. An orthant (or any rotation thereof; a simplicial cone) is not the only selfdual polyhedral cone in three or more dimensions; [21, §2.A.21] e.g, consider an *equilateral* having five extreme directions. In fact, every selfdual polyhedral cone in  $\mathbb{R}^3$  has an odd number of extreme directions. [23, thm.3]

given  $A_j \in \mathbb{S}^n$ ,  $j=1 \dots m$

$$\begin{aligned}\mathcal{K} &= \left\{ \begin{bmatrix} \langle A_1, X \rangle \\ \vdots \\ \langle A_m, X \rangle \end{bmatrix} \mid X \succeq 0 \right\} \subseteq \mathbb{R}^m \\ &= \left\{ \begin{bmatrix} \text{svec}(A_1)^T \\ \vdots \\ \text{svec}(A_m)^T \end{bmatrix} \text{svec } X \mid X \succeq 0 \right\} \\ &\triangleq \{A \text{ svec } X \mid X \succeq 0\}\end{aligned}\tag{386}$$

where  $A \in \mathbb{R}^{m \times n(n+1)/2}$ , and where symmetric vectorization  $\text{svec}$  is defined in (57). Cone  $\mathcal{K}$  is indeed convex because, by (178)

$$A \text{ svec } X_{p_1}, A \text{ svec } X_{p_2} \in \mathcal{K} \Rightarrow A(\zeta \text{ svec } X_{p_1} + \xi \text{ svec } X_{p_2}) \in \mathcal{K} \text{ for all } \zeta, \xi \geq 0 \tag{387}$$

since a nonnegatively weighted sum of positive semidefinite matrices must be positive semidefinite. (§A.3.1.0.2) Although matrix  $A$  is finite-dimensional,  $\mathcal{K}$  is generally not a polyhedral cone (unless  $m=1$  or  $2$ ) simply because  $X \in \mathbb{S}_+^n$ .

**Theorem.** *Inverse image closedness.* [225, prop.A.2.1.12] [343, thm.6.7]  
Given affine operator  $g : \mathbb{R}^m \rightarrow \mathbb{R}^p$ , convex set  $\mathcal{D} \subseteq \mathbb{R}^m$ , and convex set  $\mathcal{C} \subseteq \mathbb{R}^p \ni g^{-1}(\text{rel intr } \mathcal{C}) \neq \emptyset$ , then

$$\text{rel intr } g(\mathcal{D}) = g(\text{rel intr } \mathcal{D}), \quad \text{rel intr } g^{-1}\mathcal{C} = g^{-1}(\text{rel intr } \mathcal{C}), \quad \overline{g^{-1}\mathcal{C}} = g^{-1}\overline{\mathcal{C}} \tag{388}$$

◊

By this theorem, relative interior of  $\mathcal{K}$  may always be expressed

$$\text{rel intr } \mathcal{K} = \{A \text{ svec } X \mid X \succ 0\} \tag{389}$$

Because  $\dim(\text{aff } \mathcal{K}) = \dim(A \text{ svec } \mathbb{S}^n)$  (129) then, provided the vectorized  $A_j$  matrices are linearly independent,

$$\text{rel intr } \mathcal{K} = \text{intr } \mathcal{K} \tag{14}$$

meaning, cone  $\mathcal{K}$  is full-dimensional  $\Rightarrow$  dual cone  $\mathcal{K}^*$  is pointed by (314). Convex cone  $\mathcal{K}$  can be closed, by this corollary:

**Corollary.** *Cone closedness invariance.* [59, §3] [60, §3]

Given linear operator  $A : \mathbb{R}^p \rightarrow \mathbb{R}^m$  and closed convex cone  $\mathcal{X} \subseteq \mathbb{R}^p$ , convex cone  $\mathcal{K} = A(\mathcal{X})$  is closed ( $\overline{A(\mathcal{X})} = A(\mathcal{X})$ ) if and only if

$$\mathcal{N}(A) \cap \mathcal{X} = \{\mathbf{0}\} \quad \text{or} \quad \mathcal{N}(A) \cap \mathcal{X} \not\subseteq \text{rel } \partial \mathcal{X} \tag{390}$$

Otherwise,  $\overline{\mathcal{K}} = \overline{A(\mathcal{X})} \supseteq A(\overline{\mathcal{X}}) \supseteq A(\mathcal{X})$ . [343, thm.6.6]

◊

If matrix  $A$  has no nontrivial nullspace, then  $A \text{ svec } X$  is an isomorphism in  $X$  between cone  $\mathbb{S}_+^n$  and range  $\mathcal{R}(A)$  of matrix  $A$ ; (§2.2.1.0.1, §2.10.1.1) sufficient for convex cone  $\mathcal{K}$  to be closed and have relative boundary

$$\text{rel } \partial \mathcal{K} = \{A \text{ svec } X \mid X \succeq 0, X \not\succ 0\} \tag{391}$$

Now consider the (closed convex) dual cone:

$$\begin{aligned}
\mathcal{K}^* &= \{y \mid \langle z, y \rangle \geq 0 \text{ for all } z \in \mathcal{K}\} \subseteq \mathbb{R}^m \\
&= \{y \mid \langle z, y \rangle \geq 0 \text{ for all } z = A \text{ svec } X, X \succeq 0\} \\
&= \{y \mid \langle A \text{ svec } X, y \rangle \geq 0 \text{ for all } X \succeq 0\} \\
&= \{y \mid \langle \text{svec } X, A^T y \rangle \geq 0 \text{ for all } X \succeq 0\} \\
&= \{y \mid \text{svec}^{-1}(A^T y) \succeq 0\}
\end{aligned} \tag{392}$$

that follows from (384) and leads to an equally peculiar halfspace-description

$$\mathcal{K}^* = \{y \in \mathbb{R}^m \mid \sum_{j=1}^m y_j A_j \succeq 0\} \tag{393}$$

The summation inequality with respect to positive semidefinite cone  $\mathbb{S}_+^n$  is known as *linear matrix inequality*. [63] [172] [293] [403] Although we already know that the dual cone is convex (§2.13.1), *inverse image theorem* 2.1.9.0.1 certifies convexity of  $\mathcal{K}^*$  which is the inverse image of positive semidefinite cone  $\mathbb{S}_+^n$  under linear transformation  $g(y) \triangleq \sum y_j A_j$ . And although we already know that the dual cone is closed, it is certified by (388). By the *inverse image closedness theorem*, dual cone relative interior may always be expressed

$$\text{rel intr } \mathcal{K}^* = \{y \in \mathbb{R}^m \mid \sum_{j=1}^m y_j A_j \succ 0\} \tag{394}$$

Function  $g(y)$  on  $\mathbb{R}^m$  is an isomorphism when the vectorized  $A_j$  matrices are linearly independent; hence, uniquely invertible. Inverse image  $\mathcal{K}^*$  must therefore have dimension equal to  $\dim(\mathcal{R}(A^T) \cap \text{svec } \mathbb{S}_+^n)$  (50) and relative boundary

$$\text{rel } \partial \mathcal{K}^* = \{y \in \mathbb{R}^m \mid \sum_{j=1}^m y_j A_j \succeq 0, \sum_{j=1}^m y_j A_j \not\succeq 0\} \tag{395}$$

When this dimension equals  $m$ , then dual cone  $\mathcal{K}^*$  is full-dimensional

$$\text{rel intr } \mathcal{K}^* = \text{intr } \mathcal{K}^* \tag{14}$$

which implies: closure of convex cone  $\mathcal{K}$  is pointed (314).  $\square$

### 2.13.7 Dual of pointed polyhedral cone

In a subspace of  $\mathbb{R}^n$ , now we consider a pointed polyhedral cone  $\mathcal{K}$  given in terms of its extreme directions  $\Gamma_i$  arranged columnar in

$$X = [\Gamma_1 \ \Gamma_2 \ \cdots \ \Gamma_N] \in \mathbb{R}^{n \times N} \tag{282}$$

The *extremes theorem* (§2.8.1.1.1) provides the vertex-description of a pointed polyhedral cone in terms of its finite number of extreme directions and its lone vertex at the origin:

**2.13.7.0.1 Definition.** *Pointed polyhedral cone, vertex-description.*

Given pointed polyhedral cone  $\mathcal{K}$  in a subspace of  $\mathbb{R}^n$ , denoting its  $i^{\text{th}}$  extreme direction by  $\Gamma_i \in \mathbb{R}^n$  arranged in a matrix  $X$  as in (282), then that cone may be described: (87) (*confer*(191) (296))

$$\begin{aligned}
\mathcal{K} &= \{[\mathbf{0} \ X] a \zeta \mid a^T \mathbf{1} = 1, a \succeq 0, \zeta \geq 0\} \\
&= \{X a \zeta \mid a^T \mathbf{1} \leq 1, a \succeq 0, \zeta \geq 0\} \\
&= \{X b \mid b \succeq 0\} \subseteq \mathbb{R}^n
\end{aligned} \tag{396}$$

that is simply a conic hull (like (105)) of a finite number  $N$  of directions. Relative interior may always be expressed

$$\text{rel intr } \mathcal{K} = \{Xb \mid b \succ 0\} \subset \mathbb{R}^n \quad (397)$$

although  $Xb \in \text{rel intr } \mathcal{K} \not\Rightarrow b \succ 0$  unless matrix  $X$  represents a bijection onto its range. But identifying the cone's relative boundary in this manner

$$\text{rel } \partial\mathcal{K} = \{Xb \mid b \succeq 0, b \not\succ 0\} \quad (398)$$

holds only when  $X$  represents a bijection; in other words, some coefficients meeting lower bound zero ( $b \in \partial\mathbb{R}_+^N$ ) do not necessarily provide membership to the relative boundary of cone  $\mathcal{K}$ .  $\triangle$

Whenever cone  $\mathcal{K}$  is pointed, closed, and convex (not only polyhedral), then dual cone  $\mathcal{K}^*$  has a halfspace-description in terms of the extreme directions  $\Gamma_i$  of  $\mathcal{K}$ :

$$\mathcal{K}^* = \{y \mid \gamma^T y \geq 0 \text{ for all } \gamma \in \{\Gamma_i, i=1 \dots N\} \subseteq \text{rel } \partial\mathcal{K}\} \quad (399)$$

because when  $\{\Gamma_i\}$  constitutes any set of generators for  $\mathcal{K}$ , the discretization result in §2.13.4.1 allows relaxation of the requirement  $\forall x \in \mathcal{K}$  in (300) to  $\forall \gamma \in \{\Gamma_i\}$  directly.<sup>2.75</sup> That dual cone so defined is unique, identical to (300), polyhedral whenever the number of generators  $N$  is finite

$$\mathcal{K}^* = \{y \mid X^T y \succeq 0\} \subseteq \mathbb{R}^n \quad (366)$$

and is full-dimensional because  $\mathcal{K}$  is assumed pointed. But  $\mathcal{K}^*$  is not necessarily pointed unless  $\mathcal{K}$  is full-dimensional. (§2.13.1.2)

### 2.13.7.1 Facet normal & extreme direction

We see from (366) that the conically independent generators of cone  $\mathcal{K}$  (namely, the extreme directions of pointed closed convex cone  $\mathcal{K}$  constituting the  $N$  columns of  $X$ ) each define an inward-normal to a hyperplane supporting dual cone  $\mathcal{K}^*$  (§2.4.2.6.1) and exposing a dual facet when  $N$  is finite. Were  $\mathcal{K}^*$  pointed and finitely generated, then by closure the conjugate statement would also hold; *id est*, the extreme directions of pointed  $\mathcal{K}^*$  each define an inward-normal to a hyperplane supporting  $\mathcal{K}$  and exposing a facet when  $N$  is finite. Examine Figure 60 or Figure 66, for example.

We may conclude, the extreme directions of proper polyhedral  $\mathcal{K}$  are respectively orthogonal to the facets of  $\mathcal{K}^*$ ; likewise, the extreme directions of proper polyhedral  $\mathcal{K}^*$  are respectively orthogonal to the facets of  $\mathcal{K}$ .

## 2.13.8 Biorthogonal expansion by example

### 2.13.8.0.1 Example. Relationship to dual polyhedral cone.

Simplicial cone  $\mathcal{K}$  illustrated in Figure 67 induces a partial order on  $\mathbb{R}^2$ . All points greater than  $x$  with respect to  $\mathcal{K}$ , for example, are contained in the translated cone  $x + \mathcal{K}$ . The extreme directions  $\Gamma_1$  and  $\Gamma_2$  of  $\mathcal{K}$  do not make an orthogonal set; neither do extreme directions  $\Gamma_3$  and  $\Gamma_4$  of dual cone  $\mathcal{K}^*$ ; rather, we have the *biorthogonality condition* [406]

$$\begin{aligned} \Gamma_4^T \Gamma_1 &= \Gamma_3^T \Gamma_2 = 0 \\ \Gamma_3^T \Gamma_1 &\neq 0, \quad \Gamma_4^T \Gamma_2 \neq 0 \end{aligned} \quad (400)$$

---

<sup>2.75</sup>The extreme directions of  $\mathcal{K}$  constitute a minimal set of generators. Formulae and conversions to vertex-description of the dual cone are in §2.13.10 and §2.13.12.

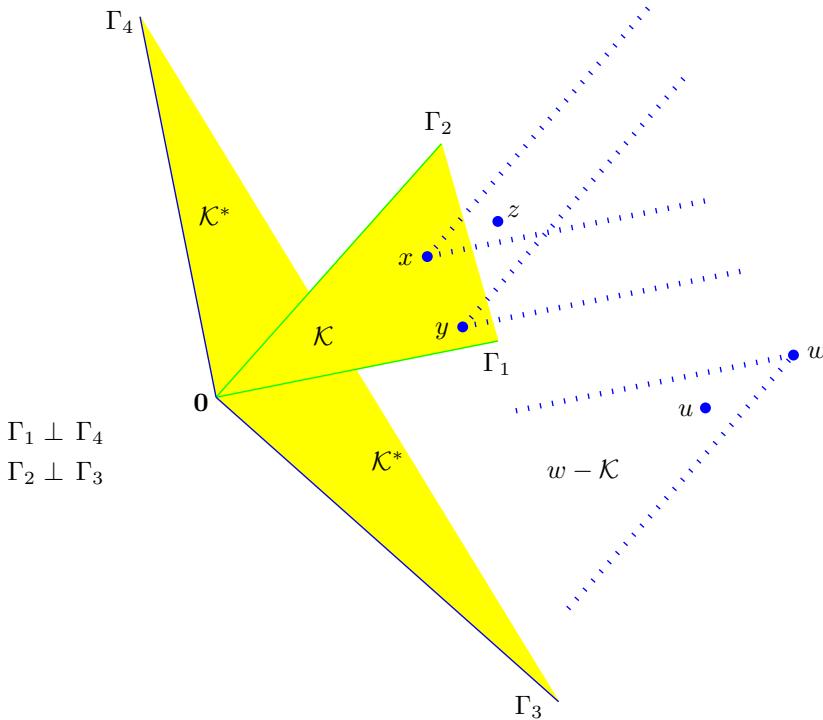


Figure 67: (confer Figure 189) Simplicial cone  $\mathcal{K} \in \mathbb{R}^2$  and its dual  $\mathcal{K}^*$  drawn truncated. Conically independent generators  $\Gamma_1$  and  $\Gamma_2$  constitute extreme directions of  $\mathcal{K}$  while  $\Gamma_3$  and  $\Gamma_4$  constitute extreme directions of  $\mathcal{K}^*$ . Dotted ray-pairs bound translated cones  $\mathcal{K}$ . Point  $x$  is comparable to point  $z$  (and *vice versa*) but not to  $y$ ;  $z \succeq_{\mathcal{K}} x \Leftrightarrow z - x \in \mathcal{K} \Leftrightarrow z - x \succeq_{\mathcal{K}} 0$  iff  $\exists$  nonnegative coordinates for biorthogonal expansion of  $z - x$ . Point  $y$  is not comparable to  $z$  because  $z$  does not belong to  $y \pm \mathcal{K}$ . Translating a negated cone is quite helpful for visualization:  $u \preceq_{\mathcal{K}} w \Leftrightarrow u \in w - \mathcal{K} \Leftrightarrow u - w \preceq_{\mathcal{K}} 0$ . Points need not belong to  $\mathcal{K}$  to be comparable; *e.g.*, all points less than  $w$  (w.r.t  $\mathcal{K}$ ) belong to  $w - \mathcal{K}$ .

Biorthogonal expansion of  $x \in \mathcal{K}$  is then

$$x = \Gamma_1 \frac{\Gamma_3^T x}{\Gamma_3^T \Gamma_1} + \Gamma_2 \frac{\Gamma_4^T x}{\Gamma_4^T \Gamma_2} \quad (401)$$

where  $\Gamma_3^T x / (\Gamma_3^T \Gamma_1)$  is the nonnegative coefficient of nonorthogonal projection (§E.6.1) of  $x$  on  $\Gamma_1$  in the direction orthogonal to  $\Gamma_3$  ( $y$  in Figure 189 p.581), and where  $\Gamma_4^T x / (\Gamma_4^T \Gamma_2)$  is the nonnegative coefficient of nonorthogonal projection of  $x$  on  $\Gamma_2$  in the direction orthogonal to  $\Gamma_4$  ( $z$  in Figure 189); they are coordinates in this nonorthogonal system. Those coefficients must be nonnegative  $x \succeq_{\mathcal{K}} 0$  because  $x \in \mathcal{K}$  (328) and  $\mathcal{K}$  is simplicial.

If we ascribe the extreme directions of  $\mathcal{K}$  to the columns of a matrix

$$X \triangleq [\Gamma_1 \ \Gamma_2] \quad (402)$$

then we find that the pseudoinverse transpose matrix

$$X^{\dagger T} = \begin{bmatrix} \Gamma_3 \frac{1}{\Gamma_3^T \Gamma_1} & \Gamma_4 \frac{1}{\Gamma_4^T \Gamma_2} \end{bmatrix} \quad (403)$$

holds the extreme directions of the dual cone. Therefore

$$x = X X^{\dagger} x \quad (409)$$

is biorthogonal expansion (401) (§E.0.1), and biorthogonality condition (400) can be expressed succinctly (§E.1.1)<sup>2.76</sup>

$$X^{\dagger} X = I \quad (410)$$

Expansion  $w = X X^{\dagger} w$ , for any particular  $w \in \mathbb{R}^n$  more generally, is unique w.r.t  $X$  if and only if the extreme directions of  $\mathcal{K}$  populating the columns of  $X \in \mathbb{R}^{n \times N}$  are linearly independent; *id est*, iff  $X$  has no nullspace.  $\square$

▼

### 2.13.8.0.2 Exercise. Visual comparison of real sums.

Given  $y \preceq x$  with respect to the nonnegative orthant, draw a figure showing a negated shifted orthant (like the cone in Figure 67) illustrating why  $\mathbf{1}^T y \leq \mathbf{1}^T x$  for  $y$  and  $x$  anywhere in  $\mathbb{R}^2$ . Incorporate two hyperplanes into your drawing, one containing  $y$  and another containing  $x$  with reference to Figure 29. Does this result hold in higher dimension?

▼

#### 2.13.8.1 Pointed cones and biorthogonality

Biorthogonality condition  $X^{\dagger} X = I$  from Example 2.13.8.0.1 means  $\Gamma_1$  and  $\Gamma_2$  are linearly independent generators of  $\mathcal{K}$  (§B.1.1.1); generators because every  $x \in \mathcal{K}$  is their conic combination. From §2.10.2 we know that means  $\Gamma_1$  and  $\Gamma_2$  must be extreme directions of  $\mathcal{K}$ .

A biorthogonal expansion is necessarily associated with a pointed closed convex cone; pointed, otherwise there can be no extreme directions (§2.8.1). We will address biorthogonal expansion with respect to a pointed polyhedral cone, not full-dimensional, in §2.13.9.

---

<sup>2.76</sup>Possibly confusing is the fact that formula  $XX^{\dagger}x$  is simultaneously: the orthogonal projection of  $x$  on  $\mathcal{R}(X)$  (2094), and a sum of nonorthogonal projections of  $x \in \mathcal{R}(X)$  on the range of each and every column of full-rank  $X$  thin-or-square (§E.5.0.0.2).

**2.13.8.1.1 Example.** *Expansions implied by diagonalization.* (confer §6.4.3.2.1)  
When matrix  $X \in \mathbb{R}^{M \times M}$  is diagonalizable (§A.5),

$$X = S\Lambda S^{-1} = [s_1 \cdots s_M] \Lambda \begin{bmatrix} w_1^T \\ \vdots \\ w_M^T \end{bmatrix} = \sum_{i=1}^M \lambda_i s_i w_i^T \quad (1699)$$

coordinates for biorthogonal expansion are its eigenvalues  $\lambda_i$  (contained in diagonal matrix  $\Lambda$ ) when expanded in  $S$ ;

$$X = SS^{-1}X = [s_1 \cdots s_M] \begin{bmatrix} w_1^T X \\ \vdots \\ w_M^T X \end{bmatrix} = \sum_{i=1}^M \lambda_i s_i w_i^T \quad (404)$$

Coordinate values depend upon geometric relationship of  $X$  to its linearly independent eigenmatrices  $s_i w_i^T$ . (§A.5.0.3, §B.1.1)

- Eigenmatrices  $s_i w_i^T$  are linearly independent dyads constituted by right and left eigenvectors of diagonalizable  $X$  and are generators of some pointed polyhedral cone  $\mathcal{K}$  in a subspace of  $\mathbb{R}^{M \times M}$ .

When  $S$  is real and  $X$  belongs to that polyhedral cone  $\mathcal{K}$ , for example, then coordinates of expansion (the eigenvalues  $\lambda_i$ ) must be nonnegative.

When matrix  $X = Q\Lambda Q^T$  is symmetric, it is diagonalizable (§A.5.1). Coordinates for biorthogonal expansion are its eigenvalues when expanded in  $Q$ ; *id est*, for  $X \in \mathbb{S}^M$

$$X = QQ^TX = \sum_{i=1}^M q_i q_i^T X = \sum_{i=1}^M \lambda_i q_i q_i^T \in \mathbb{S}^M \quad (405)$$

becomes an orthogonal expansion with *orthonormality condition*  $Q^T Q = I$  where  $\lambda_i$  is the  $i^{\text{th}}$  (largest, usually) eigenvalue of  $X$ ,  $q_i$  is the corresponding  $i^{\text{th}}$  eigenvector arranged columnar in orthogonal matrix

$$Q = [q_1 \ q_2 \ \cdots \ q_M] \in \mathbb{R}^{M \times M} \quad (406)$$

and where eigenmatrix  $q_i q_i^T$  is an extreme direction of some pointed polyhedral cone  $\mathcal{K} \subset \mathbb{S}^M$  and an extreme direction of the positive semidefinite cone  $\mathbb{S}_+^M$ .

- Orthogonal expansion is a special case of biorthogonal expansion of  $X \in \text{aff } \mathcal{K}$  occurring when polyhedral cone  $\mathcal{K}$  is any rotation about the origin of an orthant belonging to a subspace.

Similarly, when  $X = Q\Lambda Q^T$  belongs to the positive semidefinite cone in the subspace of symmetric matrices, coordinates for orthogonal expansion must be its nonnegative eigenvalues (1604) when expanded in  $Q$ ; *id est*, for  $X \in \mathbb{S}_+^M$

$$X = QQ^TX = \sum_{i=1}^M q_i q_i^T X = \sum_{i=1}^M \lambda_i q_i q_i^T \in \mathbb{S}_+^M \quad (407)$$

where  $\lambda_i \geq 0$  is the  $i^{\text{th}}$  eigenvalue of  $X$ . This means matrix  $X$  simultaneously belongs to the positive semidefinite cone and to the pointed polyhedral cone  $\mathcal{K}$  formed by the conic hull of its eigenmatrices.  $\square$

### 2.13.8.1.2 Example. Expansion respecting nonpositive orthant.

Suppose  $x \in \mathcal{K}$  any orthant in  $\mathbb{R}^n$ .<sup>2.77</sup> Then coordinates for biorthogonal expansion of  $x$  must be nonnegative; in fact, absolute value of the Cartesian coordinates.

Suppose, in particular,  $x$  belongs to the nonpositive orthant  $\mathcal{K} = \mathbb{R}_-^n$ . Then biorthogonal expansion becomes orthogonal expansion

$$x = XX^T x = \sum_{i=1}^n -e_i(-e_i^T x) = \sum_{i=1}^n -e_i |e_i^T x| \in \mathbb{R}_-^n \quad (408)$$

and the coordinates of expansion are nonnegative. For this orthant  $\mathcal{K}$  we have orthonormality condition  $X^T X = I$  where  $X = -I$ ,  $e_i \in \mathbb{R}^n$  is a standard basis vector, and  $-e_i$  is an extreme direction ([§2.8.1](#)) of  $\mathcal{K}$ .

Of course, this expansion  $x = XX^T x$  applies more broadly to domain  $\mathbb{R}^n$ , but then the coordinates each belong to all of  $\mathbb{R}$ .  $\square$

### 2.13.9 Biorthogonal expansion, derivation

Biorthogonal expansion is a means for determining coordinates in a pointed conic coordinate system characterized by a nonorthogonal basis. Study of nonorthogonal bases invokes pointed polyhedral cones and their duals; extreme directions of a cone  $\mathcal{K}$  are assumed to constitute the *basis* while those of the dual cone  $\mathcal{K}^*$  determine coordinates.

Unique biorthogonal expansion with respect to  $\mathcal{K}$  relies upon existence of its linearly independent extreme directions: Polyhedral cone  $\mathcal{K}$  must be pointed; then it possesses extreme directions. Those extreme directions must be linearly independent to uniquely represent any point in their span.

We consider nonempty pointed polyhedral cone  $\mathcal{K}$  possibly not full-dimensional; *id est*, we consider a basis spanning a subspace. Then we need only observe that section of dual cone  $\mathcal{K}^*$  in the affine hull of  $\mathcal{K}$  because, by *expansion* of  $x$ , membership  $x \in \text{aff } \mathcal{K}$  is implicit and because any breach of the ordinary dual cone into ambient space becomes irrelevant ([§2.13.10.3](#)). *Biorthogonal expansion*

$$x = XX^\dagger x \in \text{aff } \mathcal{K} = \text{aff cone}(X) \quad (409)$$

is expressed in the extreme directions  $\{\Gamma_i\}$  of  $\mathcal{K}$  arranged columnar in

$$X = [\Gamma_1 \ \Gamma_2 \ \cdots \ \Gamma_N] \in \mathbb{R}^{n \times N} \quad (282)$$

under assumption of *biorthogonality*

$$X^\dagger X = I \quad (410)$$

where  $\dagger$  denotes matrix pseudoinverse ([§E](#)).

We therefore seek, in this section, a vertex-description for  $\mathcal{K}^* \cap \text{aff } \mathcal{K}$  in terms of linearly independent dual generators  $\{\Gamma_i^*\} \subset \text{aff } \mathcal{K}$  in the same finite quantity<sup>2.78</sup> as the extreme directions  $\{\Gamma_i\}$  of

$$\mathcal{K} = \text{cone}(X) = \{Xa \mid a \succeq 0\} \subseteq \mathbb{R}^n \quad (105)$$

We assume the quantity of extreme directions  $N$  does not exceed the dimension  $n$  of ambient vector space because, otherwise, expansion w.r.t  $\mathcal{K}$  could not be unique; *id est*, assume  $N$  linearly independent extreme directions hence  $N \leq n$  ( $X$  [thin<sup>2.79</sup>](#)-or-square full-rank). In other words, wide full-rank matrix  $X$  is prohibited by uniqueness because of existence of an infinity of right inverses;

<sup>2.77</sup>An orthant is simplicial and selfdual.

<sup>2.78</sup>When  $\mathcal{K}$  is contained in a proper subspace of  $\mathbb{R}^n$ , the ordinary dual cone  $\mathcal{K}^*$  will have more generators in any minimal set than  $\mathcal{K}$  has extreme directions.

<sup>2.79</sup>“Thin” meaning more rows than columns.

- polyhedral cones whose extreme directions number in excess of the ambient space dimension are precluded in biorthogonal expansion.

### 2.13.9.1 $x \in \mathcal{K}$

Suppose  $x$  belongs to  $\mathcal{K} \subseteq \mathbb{R}^n$ . Then  $x = Xa$  for some  $a \succeq 0$ . Coordinate vector  $a$  is unique only when  $\{\Gamma_i\}$  is a linearly independent set.<sup>2.80</sup> Vector  $a \in \mathbb{R}^N$  can take the form  $a = Bx$  if  $\mathcal{R}(B) = \mathbb{R}^N$ . Then we require  $Xa = XBx = x$  and  $Bx = BXa = a$ . The pseudoinverse  $B = X^\dagger \in \mathbb{R}^{N \times n}$  ([§E](#)) is suitable when  $X$  is thin-or-square and full-rank. In that case  $\text{rank } X = N$ , and for all  $c \succeq 0$  and  $i = 1 \dots N$

$$a \succeq 0 \Leftrightarrow X^\dagger Xa \succeq 0 \Leftrightarrow a^T X^T X^\dagger c \geq 0 \Leftrightarrow \Gamma_i^T X^\dagger c \geq 0 \quad (411)$$

The penultimate inequality follows from the *generalized inequality and membership corollary*, while the last inequality is a consequence of that corollary's discretization ([§2.13.4.2.1](#)).<sup>2.81</sup> From (411) and (399) we deduce

$$\mathcal{K}^* \cap \text{aff } \mathcal{K} = \text{cone}(X^\dagger c) = \{X^\dagger c \mid c \succeq 0\} \subseteq \mathbb{R}^n \quad (412)$$

is the vertex-description for that section of  $\mathcal{K}^*$  in the affine hull of  $\mathcal{K}$  because  $\mathcal{R}(X^\dagger) = \mathcal{R}(X)$  by definition of the pseudoinverse. From (314), we know  $\mathcal{K}^* \cap \text{aff } \mathcal{K}$  must be pointed if  $\text{relintr } \mathcal{K}$  is logically assumed nonempty with respect to  $\text{aff } \mathcal{K}$ .

Conversely, suppose full-rank thin-or-square matrix ( $N \leq n$ )

$$X^{\dagger T} \triangleq [\Gamma_1^* \ \Gamma_2^* \ \dots \ \Gamma_N^*] \in \mathbb{R}^{n \times N} \quad (413)$$

comprises the extreme directions  $\{\Gamma_i^*\} \subset \text{aff } \mathcal{K}$  of the dual-cone intersection with the affine hull of  $\mathcal{K}$ .<sup>2.82</sup> From the *discretized membership theorem* and (319) we get a partial dual to (399); *id est*, assuming  $x \in \text{aff cone } X$

$$x \in \mathcal{K} \Leftrightarrow \gamma^*{}^T x \geq 0 \text{ for all } \gamma^* \in \{\Gamma_i^*, i = 1 \dots N\} \subset \partial \mathcal{K}^* \cap \text{aff } \mathcal{K} \quad (414)$$

$$\Leftrightarrow X^\dagger x \succeq 0 \quad (415)$$

that leads to a partial halfspace-description,

$$\mathcal{K} = \{x \in \text{aff cone } X \mid X^\dagger x \succeq 0\} \quad (416)$$

For  $\gamma^* = X^{\dagger T} e_i$ , any  $x = Xa$ , and for all  $i$  we have  $e_i^T X^\dagger Xa = e_i^T a \geq 0$  only when  $a \succeq 0$ . Hence  $x \in \mathcal{K}$ .

When  $X$  is full-rank, then unique biorthogonal expansion of  $x \in \mathcal{K}$  becomes (409)

$$x = XX^\dagger x = \sum_{i=1}^N \Gamma_i \Gamma_i^* {}^T x \quad (417)$$

<sup>2.80</sup> Conic independence alone ([§2.10](#)) is insufficient to guarantee uniqueness.

<sup>2.81</sup>  $a \succeq 0 \Leftrightarrow a^T X^T X^\dagger c \geq 0 \quad \forall (c \succeq 0 \Leftrightarrow a^T X^T X^\dagger c \geq 0 \quad \forall a \succeq 0)$

$$\forall (c \succeq 0 \Leftrightarrow \Gamma_i^T X^\dagger c \geq 0 \quad \forall i) \quad \blacklozenge$$

Intuitively, any nonnegative vector  $a$  is a conic combination of the standard basis  $\{e_i \in \mathbb{R}^N\}$ ;  $a \succeq 0 \Leftrightarrow a_i e_i \succeq 0$  for all  $i$ . The last inequality in (411) is a consequence of the fact that  $x = Xa$  may be any extreme direction of  $\mathcal{K}$ , in which case  $a$  is a standard basis vector;  $a = e_i \succeq 0$ . Theoretically, because  $c \succeq 0$  defines a pointed polyhedral cone (in fact, the nonnegative orthant in  $\mathbb{R}^N$ ), we can take (411) one step further by discretizing  $c$ :

$$a \succeq 0 \Leftrightarrow \Gamma_i^T \Gamma_j^* \geq 0 \text{ for } i, j = 1 \dots N \Leftrightarrow X^\dagger X \geq 0$$

In words,  $X^\dagger X$  must be a matrix whose entries are each nonnegative.

<sup>2.82</sup> When closed convex cone  $\mathcal{K}$  is not full-dimensional,  $\mathcal{K}^*$  has no extreme directions. ([§2.13.1.2](#))

whose *coordinates*  $a = \Gamma_i^{*T} x$  must be nonnegative because  $\mathcal{K}$  is assumed pointed, closed, and convex. Whenever  $X$  is full-rank, so is its pseudoinverse  $X^\dagger$ . (§E) In the present case, the columns of  $X^{\dagger T}$  are linearly independent and generators of the dual cone  $\mathcal{K}^* \cap \text{aff } \mathcal{K}$ ; hence, the columns constitute its extreme directions. (§2.10.2) That section of the dual cone is itself a polyhedral cone (by (289) or the *cone intersection theorem*, §2.7.2.1.1) having the same number of extreme directions as  $\mathcal{K}$ .

### 2.13.9.2 $x \in \text{aff } \mathcal{K}$

The extreme directions of  $\mathcal{K}$  and  $\mathcal{K}^* \cap \text{aff } \mathcal{K}$  have a distinct relationship; because  $X^\dagger X = I$ , then for  $i, j = 1 \dots N$ ,  $\Gamma_i^T \Gamma_j^* = 1$ , while for  $i \neq j$ ,  $\Gamma_i^T \Gamma_j^* = 0$ . Yet neither set of extreme directions,  $\{\Gamma_i\}$  nor  $\{\Gamma_i^*\}$ , is necessarily orthogonal. This is a biorthogonality condition, precisely, [406, §2.2.4] [228] implying each set of extreme directions is linearly independent. (§B.1.1.1)

Biorthogonal expansion therefore applies more broadly; meaning, for any  $x \in \text{aff } \mathcal{K}$ , vector  $x$  can be uniquely expressed  $x = Xb$  where  $b \in \mathbb{R}^N$  because  $\text{aff } \mathcal{K}$  contains the origin. Thus, for any such  $x \in \mathcal{R}(X)$  (confer §E.1.1), biorthogonal expansion (417) becomes  $x = XX^\dagger Xb = Xb$ .

## 2.13.10 Formulae finding dual cone

### 2.13.10.1 Pointed $\mathcal{K}$ , dual, $X$ thin-or-square full-rank

We wish to derive expressions for a convex cone and its ordinary dual under the general assumptions: pointed polyhedral  $\mathcal{K}$  denoted by its linearly independent extreme directions arranged columnar in matrix  $X$  such that

$$\text{rank}(X \in \mathbb{R}^{n \times N}) = N \triangleq \dim \text{aff } \mathcal{K} \leq n \quad (418)$$

The vertex-description is given:

$$\mathcal{K} = \{Xa \mid a \succeq 0\} \subseteq \mathbb{R}^n \quad (105)$$

from which a halfspace-description for the dual cone follows directly:

$$\mathcal{K}^* = \{y \in \mathbb{R}^n \mid X^T y \succeq 0\} \quad (419)$$

By defining a matrix

$$X^\perp \triangleq \text{basis } \mathcal{N}(X^T) \quad (420)$$

(a columnar basis for the orthogonal complement of  $\mathcal{R}(X)$ ), we can say

$$\text{aff cone } X = \text{aff } \mathcal{K} = \{x \mid X^{\perp T} x = \mathbf{0}\} \quad (421)$$

meaning  $\mathcal{K}$  lies in a subspace, perhaps  $\mathbb{R}^n$ . Thus a halfspace-description

$$\mathcal{K} = \{x \in \mathbb{R}^n \mid X^\dagger x \succeq 0, X^{\perp T} x = \mathbf{0}\} \quad (422)$$

and a vertex-description<sup>2.83</sup> from (319)

$$\mathcal{K}^* = \{[X^{\dagger T} \ X^\perp \ -X^\perp]b \mid b \succeq 0\} \subseteq \mathbb{R}^n \quad (423)$$

These results are summarized for a pointed polyhedral cone, having linearly independent generators, and its ordinary dual:

Cone Table 1	$\mathcal{K}$	$\mathcal{K}^*$
vertex-description	$X$	$X^{\dagger T}, \pm X^\perp$
halfspace-description	$X^\dagger, X^{\perp T}$	$X^T$

<sup>2.83</sup>These descriptions are not unique. A vertex-description of the dual cone, for example, might use four conically independent generators for a plane (§2.10.0.1, Figure 52) when only three would suffice.

### 2.13.10.2 Simplicial case

When a convex cone is simplicial (§2.12.3.1.1), Cone Table 1 simplifies because then  $\text{aff cone } X = \mathbb{R}^n$ : For square  $X$  and assuming simplicial  $\mathcal{K}$  such that

$$\text{rank}(X \in \mathbb{R}^{n \times N}) = N \triangleq \dim \text{aff } \mathcal{K} = n \quad (424)$$

we have

Cone Table S	$\mathcal{K}$	$\mathcal{K}^*$
vertex-description	$X$	$X^{\dagger T}$
halfspace-description	$X^\dagger$	$X^T$

For example, vertex-description (423) simplifies to

$$\mathcal{K}^* = \{X^{\dagger T} b \mid b \succeq 0\} \subset \mathbb{R}^n \quad (425)$$

Now, because  $\dim \mathcal{R}(X) = \dim \mathcal{R}(X^{\dagger T})$ , (§E) dual cone  $\mathcal{K}^*$  is simplicial whenever  $\mathcal{K}$  is. So (§2.10.2) each respective vertex-description holds the extreme directions of the corresponding cone.

### 2.13.10.3 Cone membership relations in a subspace $\mathcal{S}_{\mathcal{R}}$

It is obvious by definition (300) of ordinary dual cone  $\mathcal{K}^*$ , in ambient vector space  $\mathcal{R}$ , that its determination instead in subspace  $\mathcal{S}_{\mathcal{R}} \subseteq \mathcal{R}$  is identical to its intersection with  $\mathcal{S}_{\mathcal{R}}$ ; *id est*, assuming closed convex cone  $\mathcal{K} \subseteq \mathcal{S}_{\mathcal{R}}$  and  $\mathcal{K}^* \subseteq \mathcal{R}$

$$(\mathcal{K}^* \text{ were ambient } \mathcal{S}_{\mathcal{R}}) \equiv (\mathcal{K}^* \text{ in ambient } \mathcal{R}) \cap \mathcal{S}_{\mathcal{R}} \quad (426)$$

because

$$\{y \in \mathcal{S}_{\mathcal{R}} \mid \langle y, x \rangle \geq 0 \text{ for all } x \in \mathcal{K}\} = \{y \in \mathcal{R} \mid \langle y, x \rangle \geq 0 \text{ for all } x \in \mathcal{K}\} \cap \mathcal{S}_{\mathcal{R}} \quad (427)$$

From this, a constrained membership relation for the ordinary dual cone  $\mathcal{K}^* \subseteq \mathcal{R}$ , assuming  $x, y \in \mathcal{S}_{\mathcal{R}}$  and closed convex cone  $\mathcal{K} \subseteq \mathcal{S}_{\mathcal{R}}$

$$y \in \mathcal{K}^* \cap \mathcal{S}_{\mathcal{R}} \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } x \in \mathcal{K} \quad (428)$$

By closure in subspace  $\mathcal{S}_{\mathcal{R}}$  we have conjugation (§2.13.1.2):

$$x \in \mathcal{K} \Leftrightarrow \langle y, x \rangle \geq 0 \text{ for all } y \in \mathcal{K}^* \cap \mathcal{S}_{\mathcal{R}} \quad (429)$$

This means membership determination in subspace  $\mathcal{S}_{\mathcal{R}}$  requires knowledge of dual cone only in  $\mathcal{S}_{\mathcal{R}}$ . For sake of completeness, for proper cone  $\mathcal{K}$  with respect to subspace  $\mathcal{S}_{\mathcal{R}}$  (confer (329))

$$x \in \text{intr } \mathcal{K} \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } y \in \mathcal{K}^* \cap \mathcal{S}_{\mathcal{R}}, y \neq \mathbf{0} \quad (430)$$

$$x \in \mathcal{K}, x \neq \mathbf{0} \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } y \in \text{intr } \mathcal{K}^* \cap \mathcal{S}_{\mathcal{R}} \quad (431)$$

(By closure, we also have the conjugate relations.) Yet when  $\mathcal{S}_{\mathcal{R}}$  equals  $\text{aff } \mathcal{K}$  for  $\mathcal{K}$  a closed convex cone

$$x \in \text{relintr } \mathcal{K} \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } y \in \mathcal{K}^* \cap \text{aff } \mathcal{K}, y \neq \mathbf{0} \quad (432)$$

$$x \in \mathcal{K}, x \neq \mathbf{0} \Leftrightarrow \langle y, x \rangle > 0 \text{ for all } y \in \text{relintr}(\mathcal{K}^* \cap \text{aff } \mathcal{K}) \quad (433)$$

#### 2.13.10.4 Subspace $\mathcal{S}_{\mathcal{R}} = \text{aff } \mathcal{K}$

Assume now a subspace  $\mathcal{S}_{\mathcal{R}}$  that is the affine hull of cone  $\mathcal{K}$ : Consider again a pointed polyhedral cone  $\mathcal{K}$  denoted by its extreme directions arranged columnar in matrix  $X$  such that

$$\text{rank}(X \in \mathbb{R}^{n \times N}) = N \triangleq \dim \text{aff } \mathcal{K} \leq n \quad (418)$$

We want expressions for the convex cone and its dual in subspace  $\mathcal{S}_{\mathcal{R}} = \text{aff } \mathcal{K}$ :

Cone Table A		$\mathcal{K}$	$\mathcal{K}^* \cap \text{aff } \mathcal{K}$
vertex-description		$X$	$X^{\dagger T}$
halfspace-description		$X^{\dagger}, X^{\perp T}$	$X^T, X^{\perp T}$

Now each respective vertex-description holds extreme directions of the corresponding cone in subspace  $\mathcal{S}_{\mathcal{R}}$ . (§2.13.9.1) When  $\dim \text{aff } \mathcal{K} = n$ , this table reduces to Cone Table S.

These descriptions facilitate work in a proper subspace. The subspace of symmetric matrices  $\mathbb{S}^N$ , for example, often serves as ambient space. [2.84](#)

##### 2.13.10.4.1 Exercise. Conically independent columns and rows.

We suspect the number of conically independent columns (rows) of  $X$  to be the same for  $X^{\dagger T}$ , where  $\dagger$  denotes matrix pseudoinverse (§E). Prove whether it holds that the columns (rows) of  $X$  are c.i.  $\Leftrightarrow$  the columns (rows) of  $X^{\dagger T}$  are c.i. ▼

##### 2.13.10.4.2 Example. Monotone nonnegative cone. [65, exer.2.33] [394, §2]

Simplicial cone (§2.12.3.1.1)  $\mathcal{K}_{\mathcal{M}+}$  is the cone of all nonnegative vectors having their entries sorted in nonincreasing order:

$$\begin{aligned} \mathcal{K}_{\mathcal{M}+} &\triangleq \{x \mid x_1 \geq x_2 \geq \dots \geq x_n \geq 0\} \subseteq \mathbb{R}_+^n \\ &= \{x \mid (e_i - e_{i+1})^T x \geq 0, i = 1 \dots n-1, e_n^T x \geq 0\} \\ &= \{x \mid X^{\dagger T} x \succeq 0\} \end{aligned} \quad (434)$$

a halfspace-description where  $e_i$  is the  $i^{\text{th}}$  standard basis vector, and where [2.85](#)

$$X^{\dagger T} \triangleq [e_1 - e_2 \quad e_2 - e_3 \quad \dots \quad e_n] \in \mathbb{R}^{n \times n} \quad (435)$$

For any vectors  $x$  and  $y$ , simple algebra demands

$$\begin{aligned} x^T y &= \sum_{i=1}^n x_i y_i = (x_1 - x_2)y_1 + (x_2 - x_3)(y_1 + y_2) + (x_3 - x_4)(y_1 + y_2 + y_3) + \dots \\ &\quad + (x_{n-1} - x_n)(y_1 + \dots + y_{n-1}) + x_n(y_1 + \dots + y_n) \end{aligned} \quad (436)$$

Because  $x_i - x_{i+1} \geq 0 \forall i$  by assumption whenever  $x \in \mathcal{K}_{\mathcal{M}+}$ , we can employ dual generalized inequalities (326) with respect to the selfdual nonnegative orthant  $\mathbb{R}_+^n$  to find the halfspace-description of dual monotone nonnegative cone  $\mathcal{K}_{\mathcal{M}+}^*$ . We can say  $x^T y \geq 0$  for all  $X^{\dagger T} x \succeq 0$  [sic] if and only if

$$y_1 \geq 0, \quad y_1 + y_2 \geq 0, \quad \dots, \quad y_1 + y_2 + \dots + y_n \geq 0 \quad (437)$$

[2.84](#)The dual cone of positive semidefinite matrices  $\mathbb{S}_+^{N*} = \mathbb{S}_+^N$  remains in  $\mathbb{S}^N$  by convention, whereas the ordinary dual cone would venture into  $\mathbb{R}^{N \times N}$ .

[2.85](#)With  $X^{\dagger}$  in hand, we might concisely scribe the remaining vertex- and halfspace-descriptions from the tables for  $\mathcal{K}_{\mathcal{M}+}$  and its dual. Instead we use dual generalized inequalities in their derivation.

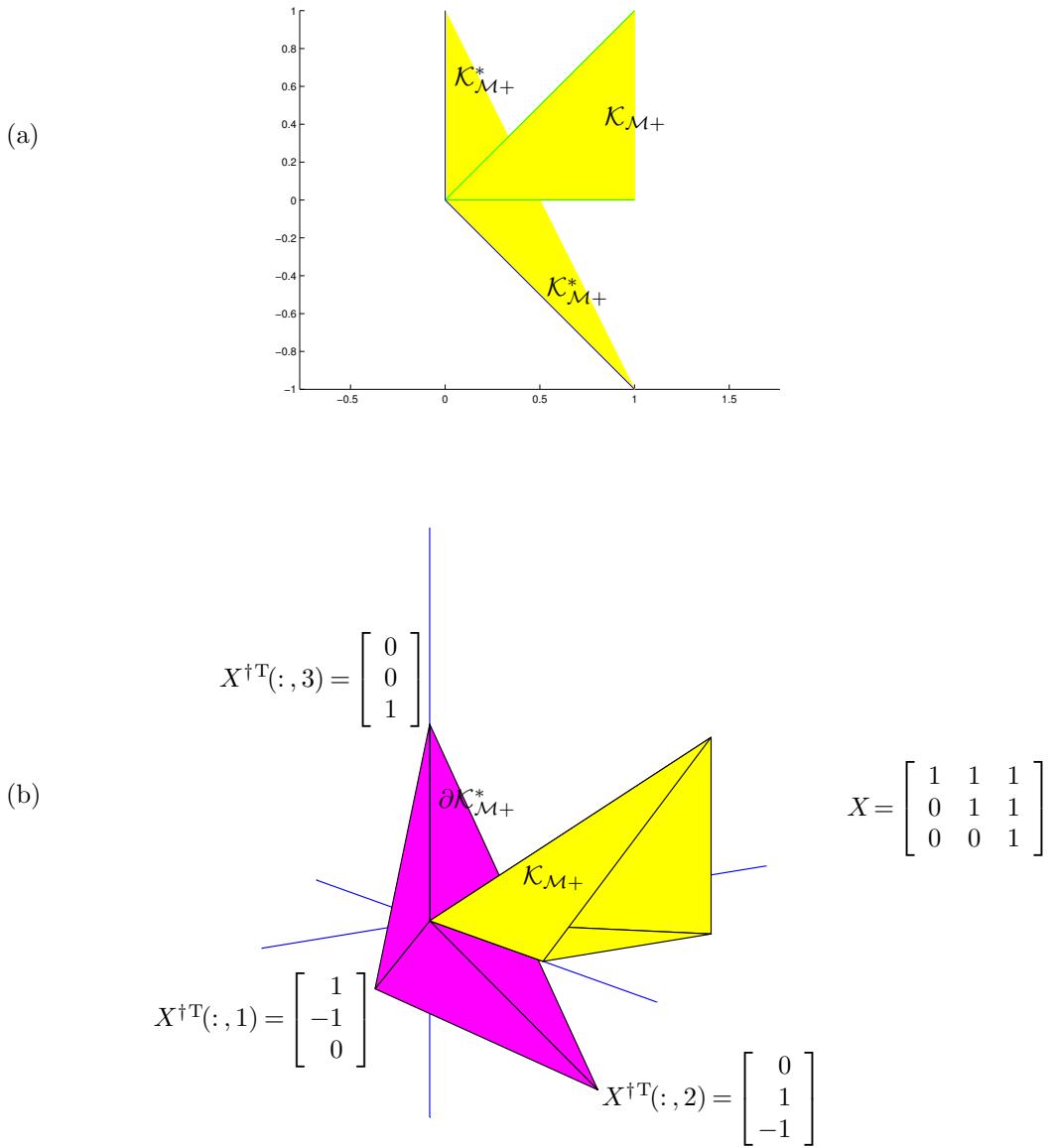
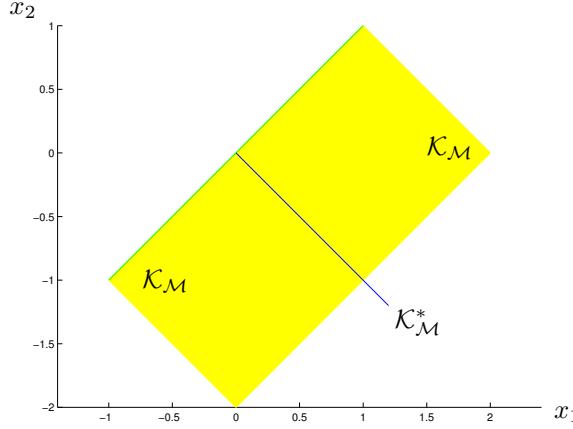


Figure 68: Simplicial cones. **(a)** Monotone nonnegative cone  $K_{\mathcal{M}+}$  and its dual  $K_{\mathcal{M}+}^*$  (drawn truncated) in  $\mathbb{R}^2$ . **(b)** Monotone nonnegative cone and boundary of its dual (both drawn truncated) in  $\mathbb{R}^3$ . Extreme directions of  $K_{\mathcal{M}+}^*$  are indicated.

Figure 69: Monotone cone  $\mathcal{K}_M$  and its dual  $\mathcal{K}_M^*$  (drawn truncated) in  $\mathbb{R}^2$ .

*id est,*

$$x^T y \geq 0 \quad \forall X^\dagger x \succeq 0 \Leftrightarrow X^T y \succeq 0 \quad (438)$$

where

$$X = [e_1 \ e_1 + e_2 \ e_1 + e_2 + e_3 \ \cdots \ \mathbf{1}] \in \mathbb{R}^{n \times n} \quad (439)$$

Because  $X^\dagger x \succeq 0$  connotes membership of  $x$  to pointed  $\mathcal{K}_{M+}$ , then (by (300)) the dual cone that we seek comprises all  $y$  for which (438) holds; thus its halfspace-description

$$\mathcal{K}_{M+}^* = \{y \succeq 0\} = \{y \mid \sum_{i=1}^k y_i \geq 0, k=1 \dots n\} = \{y \mid X^T y \succeq 0\} \subset \mathbb{R}^n \quad (440)$$

The monotone nonnegative cone and its dual are simplicial, illustrated for two Euclidean spaces in Figure 68.

From §2.13.7.1, the extreme directions of proper  $\mathcal{K}_{M+}$  are respectively orthogonal to the facets of  $\mathcal{K}_{M+}^*$ . Because  $\mathcal{K}_{M+}^*$  is simplicial, the inward-normals to its facets constitute the linearly independent rows of  $X^T$  by (440). Hence the vertex-description for  $\mathcal{K}_{M+}$  employs the columns of  $X$  in agreement with Cone Table S because  $X^\dagger = X^{-1}$ . Likewise, the extreme directions of proper  $\mathcal{K}_{M+}^*$  are respectively orthogonal to the facets of  $\mathcal{K}_{M+}$  whose inward-normals are contained in the rows of  $X^\dagger$  by (434). So the vertex-description for  $\mathcal{K}_{M+}^*$  employs the columns of  $X^{\dagger T}$ .  $\square$

#### 2.13.10.4.3 Example. Monotone cone. (Figure 69, Figure 70)

Full-dimensional but not pointed, the monotone cone is polyhedral and defined by the halfspace-description

$$\mathcal{K}_M \triangleq \{x \in \mathbb{R}^n \mid x_1 \geq x_2 \geq \cdots \geq x_n\} = \{x \in \mathbb{R}^n \mid X^{*T} x \succeq 0\} \quad (441)$$

Its dual is therefore pointed but not full-dimensional;

$$\mathcal{K}_M^* = \{X^* b \triangleq [e_1 - e_2 \ e_2 - e_3 \ \cdots \ e_{n-1} - e_n] b \mid b \succeq 0\} \subset \mathbb{R}^n \quad (442)$$

the dual cone vertex-description where the columns of  $X^*$  comprise its extreme directions. Because dual monotone cone  $\mathcal{K}_M^*$  is pointed and satisfies

$$\text{rank}(X^* \in \mathbb{R}^{n \times N}) = N \triangleq \dim \text{aff } \mathcal{K}^* \leq n \quad (443)$$

where  $N = n - 1$ , and because  $\mathcal{K}_M$  is closed and convex, we may adapt Cone Table 1 (p.151) as follows:

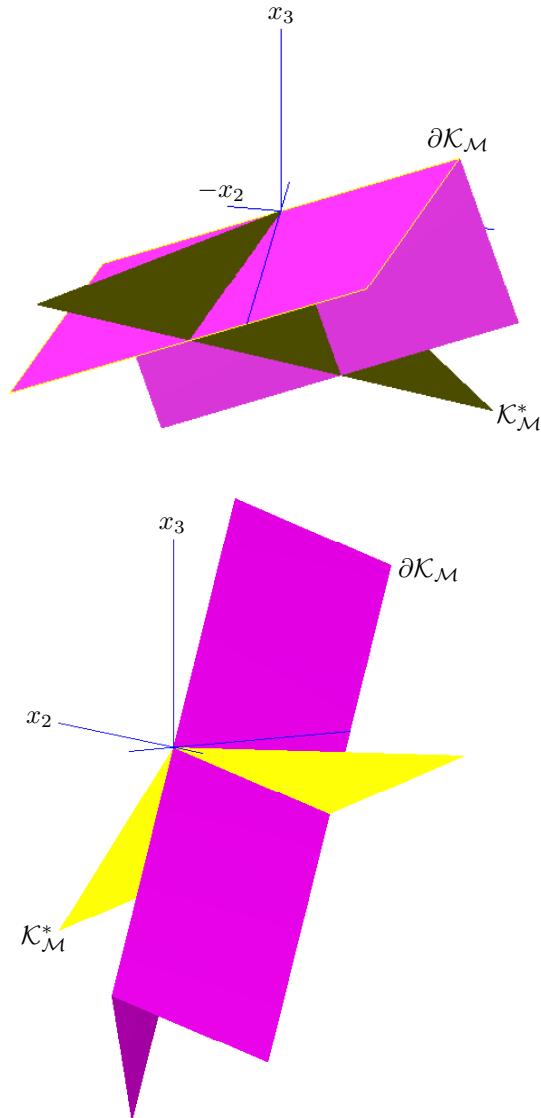


Figure 70: Two views of monotone cone  $\mathcal{K}_M$  and its dual  $\mathcal{K}_M^*$  (drawn truncated) in  $\mathbb{R}^3$ . Monotone cone is not pointed. Dual monotone cone is not full-dimensional. (Cartesian coordinate axes are drawn for reference.)

Cone Table 1*	$\mathcal{K}^*$	$\mathcal{K}^{**} = \mathcal{K}$
vertex-description	$X^*$	$X^{*\dagger T}, \pm X^{*\perp}$
halfspace-description	$X^{*\dagger}, X^{*\perp T}$	$X^{*T}$

The vertex-description for  $\mathcal{K}_M$  is therefore

$$\mathcal{K}_M = \{[X^{*\dagger T} \quad X^{*\perp} \quad -X^{*\perp}]a \mid a \succeq 0\} \subset \mathbb{R}^n \quad (444)$$

where  $X^{*\perp} = \mathbf{1}$  and

$$X^{*\dagger} = \frac{1}{n} \begin{bmatrix} n-1 & -1 & -1 & \cdots & -1 & -1 & -1 \\ n-2 & n-2 & -2 & \ddots & \cdots & -2 & -2 \\ \vdots & n-3 & n-3 & \ddots & -(n-4) & \vdots & -3 \\ 3 & \vdots & n-4 & \ddots & -(n-3) & -(n-3) & \vdots \\ 2 & 2 & \cdots & \ddots & 2 & -(n-2) & -(n-2) \\ 1 & 1 & 1 & \cdots & 1 & 1 & -(n-1) \end{bmatrix} \in \mathbb{R}^{n-1 \times n} \quad (445)$$

while

$$\mathcal{K}_M^* = \{y \in \mathbb{R}^n \mid X^{*\dagger}y \succeq 0, X^{*\perp T}y = \mathbf{0}\} \quad (446)$$

is the dual monotone cone halfspace-description.  $\square$

#### 2.13.10.4.4 Exercise. Inside the monotone cones.

Mathematically describe the respective interior of the monotone nonnegative cone and monotone cone. In three dimensions, also describe the relative interior of each face.  $\blacktriangledown$

#### 2.13.10.5 More pointed cone descriptions with equality condition

Consider pointed polyhedral cone  $\mathcal{K}$  having a linearly independent set of generators and whose subspace membership is explicit; *id est*, we are given the ordinary halfspace-description

$$\mathcal{K} = \{x \mid Ax \succeq 0, Cx = \mathbf{0}\} \subseteq \mathbb{R}^n \quad (289a)$$

where  $A \in \mathbb{R}^{m \times n}$  and  $C \in \mathbb{R}^{p \times n}$ . This can be equivalently written in terms of nullspace of  $C$  and vector  $\xi$ :

$$\mathcal{K} = \{Z\xi \in \mathbb{R}^n \mid AZ\xi \succeq 0\} \quad (447)$$

where  $\mathcal{R}(Z \in \mathbb{R}^{n \times n - \text{rank } C}) \triangleq \mathcal{N}(C)$ . Assuming (418) is satisfied

$$\text{rank } X \triangleq \text{rank}((AZ)^\dagger \in \mathbb{R}^{n - \text{rank } C \times m}) = m - \ell = \dim \text{aff } \mathcal{K} \leq n - \text{rank } C \quad (448)$$

where  $\ell$  is the number of conically dependent rows in  $AZ$  which must be removed to make  $\hat{A}Z$  before the Cone Tables become applicable.<sup>2.86</sup> Then results collected there admit assignment  $\hat{X} \triangleq (\hat{A}Z)^\dagger \in \mathbb{R}^{n - \text{rank } C \times m - \ell}$ , where  $\hat{A} \in \mathbb{R}^{m - \ell \times n}$ , followed with linear transformation by  $Z$ . So we get the vertex-description, for full-rank  $(\hat{A}Z)^\dagger$  thin-or-square,

$$\mathcal{K} = \{Z(\hat{A}Z)^\dagger b \mid b \succeq 0\} \quad (449)$$

From this and (366) we get a halfspace-description of the dual cone

$$\mathcal{K}^* = \{y \in \mathbb{R}^n \mid (Z^T \hat{A}^T)^\dagger Z^T y \succeq 0\} \quad (450)$$

---

<sup>2.86</sup>When the conically dependent rows (§2.10) are removed, the rows remaining must be linearly independent for the Cone Tables (p.12) to apply.

From this and Cone Table 1 (p.151) we get a vertex-description, (2050)

$$\mathcal{K}^* = \{[Z^{\dagger T}(\hat{A}Z)^T \quad C^T \quad -C^T]c \mid c \succeq 0\} \quad (451)$$

Yet because

$$\mathcal{K} = \{x \mid Ax \succeq 0\} \cap \{x \mid Cx = \mathbf{0}\} \quad (452)$$

then, by (319), we get an equivalent vertex-description for the dual cone

$$\begin{aligned} \mathcal{K}^* &= \overline{\{x \mid Ax \succeq 0\}^* + \{x \mid Cx = \mathbf{0}\}^*} \\ &= \{[A^T \quad C^T \quad -C^T]b \mid b \succeq 0\} \end{aligned} \quad (453)$$

from which the conically dependent columns may, of course, be removed.

### 2.13.11 Dual cone-translate

(§E.10.3.2.1) First-order optimality condition (355) inspires a dual-cone variant: For any set  $\mathcal{K}$ , the negative dual of its translation by any  $a \in \mathbb{R}^n$  is

$$\begin{aligned} -(\mathcal{K} - a)^* &= \{y \in \mathbb{R}^n \mid \langle y, x - a \rangle \leq 0 \text{ for all } x \in \mathcal{K}\} \triangleq \mathcal{K}^\perp(a) \\ &= \{y \in \mathbb{R}^n \mid \langle y, x \rangle \leq 0 \text{ for all } x \in \mathcal{K} - a\} \end{aligned} \quad (454)$$

a closed convex cone called *normal cone* to  $\mathcal{K}$  at point  $a$ . From this, a new membership relation like (323):

$$y \in -(\mathcal{K} - a)^* \Leftrightarrow \langle y, x - a \rangle \leq 0 \text{ for all } x \in \mathcal{K} \quad (455)$$

and by closure the conjugate, for closed convex cone  $\mathcal{K}$

$$x \in \mathcal{K} \Leftrightarrow \langle y, x - a \rangle \leq 0 \text{ for all } y \in -(\mathcal{K} - a)^* \quad (456)$$

#### 2.13.11.1 first-order optimality condition - restatement

(confer §2.13.3) The general first-order necessary and sufficient condition for optimality of solution  $x^*$  to a minimization problem with real differentiable convex objective function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  over convex feasible set  $\mathcal{C}$  is [342, §3]

$$-\nabla f(x^*) \in -(\mathcal{C} - x^*)^*, \quad x^* \in \mathcal{C} \quad (457)$$

*id est*, the negative gradient (§3.6) belongs to the normal cone to  $\mathcal{C}$  at  $x^*$  as in Figure 71.

##### 2.13.11.1.1 Example. Normal cone to orthant.

Consider proper cone  $\mathcal{K} = \mathbb{R}_+^n$ , the selfdual nonnegative orthant in  $\mathbb{R}^n$ . The normal cone to  $\mathbb{R}_+^n$  at  $a \in \mathcal{K}$  is (2291)

$$\mathcal{K}_{\mathbb{R}_+^n}^\perp(a \in \mathbb{R}_+^n) = -(\mathbb{R}_+^n - a)^* = -\mathbb{R}_+^n \cap a^\perp, \quad a \in \mathbb{R}_+^n \quad (458)$$

where  $-\mathbb{R}_+^n = -\mathcal{K}^*$  is the algebraic complement of  $\mathbb{R}_+^n$ , and  $a^\perp$  is the orthogonal complement to range of vector  $a$ . This means: When point  $a$  is interior to  $\mathbb{R}_+^n$ , the normal cone is the origin. If  $n_p$  represents number of nonzero entries in vector  $a \in \partial \mathbb{R}_+^n$ , then  $\dim(-\mathbb{R}_+^n \cap a^\perp) = n - n_p$  and there is a complementary relationship between the nonzero entries in vector  $a$  and the nonzero entries in any vector  $x \in -\mathbb{R}_+^n \cap a^\perp$ .  $\square$

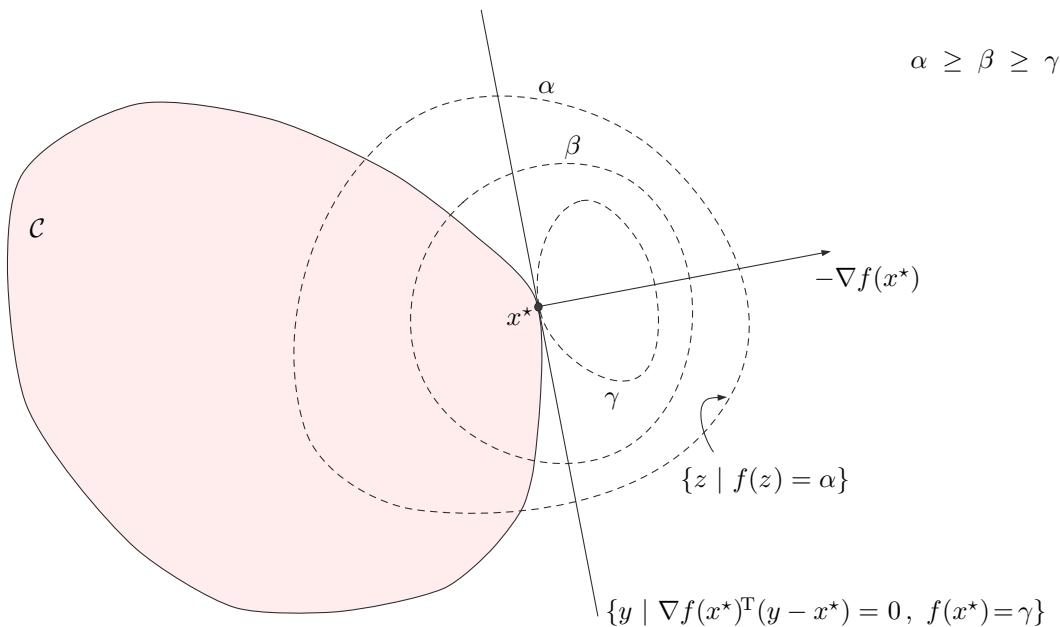


Figure 71: (confer Figure 83) Shown is a plausible contour plot in  $\mathbb{R}^2$  of some arbitrary differentiable convex real function  $f(x)$  at selected levels  $\alpha$ ,  $\beta$ , and  $\gamma$ ; *id est*, contours of equal level  $f$  (*level sets*) drawn dashed in function's domain. From results in §3.7 (p.198), gradient  $\nabla f(x^*)$  is normal to  $\gamma$ -sublevel set  $\mathcal{L}_\gamma f$  (566) by Definition E.9.1.0.1. From §2.13.11.1, function is minimized over convex set  $\mathcal{C}$  at point  $x^*$  iff negative gradient  $-\nabla f(x^*)$  belongs to normal cone to  $\mathcal{C}$  there. In circumstance depicted, normal cone is a ray whose direction is coincident with negative gradient. So, gradient is normal to a hyperplane supporting both  $\mathcal{C}$  and the  $\gamma$ -sublevel set.

**2.13.11.1.2 Example.** *Optimality conditions for conic problem.*

Consider a convex optimization problem having real differentiable convex objective function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  defined on domain  $\mathbb{R}^n$

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{subject to} && x \in \mathcal{K} \end{aligned} \quad (459)$$

Let's first suppose that the feasible set is a pointed polyhedral cone  $\mathcal{K}$  possessing a linearly independent set of generators and whose subspace membership is made explicit by wide full-rank matrix  $C \in \mathbb{R}^{p \times n}$ ; *id est*, we are given the halfspace-description, for  $A \in \mathbb{R}^{m \times n}$

$$\mathcal{K} = \{x \mid Ax \succeq 0, Cx = \mathbf{0}\} \subseteq \mathbb{R}^n \quad (289a)$$

(We'll generalize to any convex cone  $\mathcal{K}$  shortly.) Vertex-description of this cone, assuming  $(\hat{A}Z)^\dagger$  thin-or-square full-rank, is

$$\mathcal{K} = \{Z(\hat{A}Z)^\dagger b \mid b \succeq 0\} \quad (449)$$

where  $\hat{A} \in \mathbb{R}^{m-\ell \times n}$ ,  $\ell$  is the number of conically dependent rows in  $AZ$  (§2.10) which must be removed, and  $Z \in \mathbb{R}^{n \times n - \text{rank } C}$  holds basis  $\mathcal{N}(C)$  columnar.

From optimality condition (355),

$$\nabla f(x^*)^T (Z(\hat{A}Z)^\dagger b - x^*) \geq 0 \quad \forall b \succeq 0 \quad (460)$$

$$-\nabla f(x^*)^T Z(\hat{A}Z)^\dagger (b - b^*) \leq 0 \quad \forall b \succeq 0 \quad (461)$$

because

$$x^* \triangleq Z(\hat{A}Z)^\dagger b^* \in \mathcal{K} \quad (462)$$

From membership relation (455) and Example 2.13.11.1.1

$$\begin{aligned} & \langle -(Z^T \hat{A}^T)^\dagger Z^T \nabla f(x^*), b - b^* \rangle \leq 0 \quad \text{for all } b \in \mathbb{R}_+^{m-\ell} \\ & \Leftrightarrow \\ & -(Z^T \hat{A}^T)^\dagger Z^T \nabla f(x^*) \in -\mathbb{R}_+^{m-\ell} \cap b^{*\perp} \end{aligned} \quad (463)$$

Then equivalent necessary and sufficient conditions for optimality of conic problem (459) with feasible set  $\mathcal{K}$  are: (*confer*(365))

$$(Z^T \hat{A}^T)^\dagger Z^T \nabla f(x^*) \succeq_0 0, \quad \underset{\mathbb{R}_+^{m-\ell}}{b^* \succeq_0 0}, \quad \nabla f(x^*)^T Z(\hat{A}Z)^\dagger b^* = 0 \quad (464)$$

expressible, by (450),

$$\nabla f(x^*) \in \mathcal{K}^*, \quad x^* \in \mathcal{K}, \quad \nabla f(x^*)^T x^* = 0 \quad (465)$$

This result (465) actually applies more generally to any convex cone  $\mathcal{K}$  comprising the feasible set: Necessary and sufficient optimality conditions are in terms of objective gradient

$$-\nabla f(x^*) \in -(\mathcal{K} - x^*)^*, \quad x^* \in \mathcal{K} \quad (457)$$

whose membership to normal cone, assuming only cone  $\mathcal{K}$  convexity,

$$-(\mathcal{K} - x^*)^* = \mathcal{K}_{\mathcal{K}}^\perp (x^* \in \mathcal{K}) = -\mathcal{K}^* \cap x^{*\perp} \quad (2291)$$

equivalently expresses conditions (465).

When  $\mathcal{K} = \mathbb{R}_+^n$ , in particular, then  $C = \mathbf{0}$ ,  $A = Z = I \in \mathbb{S}^n$ ; *id est*,

$$\begin{array}{ll} \underset{x}{\text{minimize}} & f(x) \\ \text{subject to} & x \succeq 0 \\ & \mathbb{R}_+^n \end{array} \quad (466)$$

Necessary and sufficient optimality conditions become (*confer* [65, §4.2.3])

$$\begin{array}{lll} \nabla f(x^*) \succeq 0, & x^* \succeq 0, & \nabla f(x^*)^T x^* = 0 \\ \mathbb{R}_+^n & \mathbb{R}_+^n & \end{array} \quad (467)$$

equivalent to condition (333)<sup>2.87</sup> (under nonzero gradient) for membership to the nonnegative orthant boundary  $\partial \mathbb{R}_+^n$ .  $\square$

#### 2.13.11.1.3 Example. Complementarity problem.

[236]

A complementarity problem in nonlinear function  $f$  is nonconvex

$$\begin{array}{ll} \text{find} & z \in \mathcal{K} \\ \text{subject to} & f(z) \in \mathcal{K}^* \\ & \langle z, f(z) \rangle = 0 \end{array} \quad (468)$$

yet bears strong resemblance to (465) and to Moreau's decomposition (2228) on page 600 for projection  $P$  on mutually polar cones  $\mathcal{K}$  and  $-\mathcal{K}^*$ . Identify a sum of mutually orthogonal projections  $x \triangleq z - f(z)$ ; in Moreau's terms,  $z = P_{\mathcal{K}}x$  and  $-f(z) = P_{-\mathcal{K}^*}x$ . Then  $f(z) \in \mathcal{K}^*$  (§E.9.2.2 no.4) and  $z$  is a solution to the complementarity problem iff it is a fixed point of

$$z = P_{\mathcal{K}}x = P_{\mathcal{K}}(z - f(z)) \quad (469)$$

Given that a solution exists, existence of a fixed point would be guaranteed by theory of *contraction*. [254, p.300] But because only *nonexpansivity* (Theorem E.9.3.0.1) is achievable by a projector, uniqueness cannot be assured. [229, p.155] Elegant proofs of equivalence between complementarity problem (468) and fixed point problem (469) are provided by Németh [421, *Fixed point problems*].  $\square$

#### 2.13.11.1.4 Example. Linear complementarity problem.

[95] [303] [347]

Given matrix  $B \in \mathbb{R}^{n \times n}$  and vector  $q \in \mathbb{R}^n$ , a prototypical complementarity problem on the nonnegative orthant  $\mathcal{K} = \mathbb{R}_+^n$  is linear in  $w = f(z)$ :

$$\begin{array}{ll} \text{find} & z \succeq 0 \\ \text{subject to} & w \succeq 0 \\ & w^T z = 0 \\ & w = q + Bz \end{array} \quad (470)$$

This problem is not convex when both vectors  $w$  and  $z$  are variable.<sup>2.88</sup> Notwithstanding, this linear complementarity problem can be solved by identifying  $w \leftarrow \nabla f(z) = q + Bz$

<sup>2.87</sup> and equivalent to well-known Karush-Kuhn-Tucker (KKT) optimality conditions [65, §5.5.3] because the dual variable becomes gradient  $\nabla f(x)$ .

<sup>2.88</sup>But if one of them is fixed, then the problem becomes convex with a very simple geometric interpretation: Define the affine subset

$$\mathcal{A} \triangleq \{y \in \mathbb{R}^n \mid By = w - q\}$$

For  $w^T z$  to vanish, there must be a complementary relationship between the nonzero entries of vectors  $w$  and  $z$ ; *id est*,  $w_i z_i = 0 \forall i$ . Given  $w \succeq 0$ , then  $z$  belongs to the convex set of solutions:

$$z \in -\mathcal{K}_{\mathbb{R}_+^n}^\perp (w \in \mathbb{R}_+^n) \cap \mathcal{A} = \mathbb{R}_+^n \cap w^\perp \cap \mathcal{A}$$

where  $\mathcal{K}_{\mathbb{R}_+^n}^\perp(w)$  is the normal cone to  $\mathbb{R}_+^n$  at  $w$  (458). If this intersection is nonempty, then the problem is solvable.

then substituting that gradient into (468)

$$\begin{aligned} & \text{find } z \in \mathcal{K} \\ & \text{subject to } \nabla f(z) \in \mathcal{K}^* \\ & \quad \langle z, \nabla f(z) \rangle = 0 \end{aligned} \tag{471}$$

which is simply a restatement of optimality conditions (465) for conic problem (459). Suitable  $f(z)$  is the quadratic objective from convex problem

$$\begin{aligned} & \underset{z}{\text{minimize}} \quad \frac{1}{2} z^T B z + q^T z \\ & \text{subject to } z \succeq 0 \end{aligned} \tag{472}$$

which means  $B \in \mathbb{S}_+^n$  can be (symmetric) positive semidefinite for solution of (470) by this method. Then (470) has solution iff (472) does.  $\square$

#### 2.13.11.1.5 Exercise. Optimality for equality-constrained conic problem.

Consider a conic optimization problem like (459) having real differentiable convex objective function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$

$$\begin{aligned} & \underset{x}{\text{minimize}} \quad f(x) \\ & \text{subject to } Cx = d \\ & \quad x \in \mathcal{K} \end{aligned} \tag{473}$$

minimized over convex cone  $\mathcal{K}$  but, this time, constrained to affine set  $\mathcal{A} = \{x \mid Cx = d\}$ . Show, by means of first-order optimality condition (355) or (457), that necessary and sufficient optimality conditions are: (*confer* (465))

$$\begin{aligned} & x^* \in \mathcal{K} \\ & Cx^* = d \\ & \nabla f(x^*) + C^T \nu^* \in \mathcal{K}^* \\ & \langle \nabla f(x^*) + C^T \nu^*, x^* \rangle = 0 \end{aligned} \tag{474}$$

where  $\nu^*$  is any vector<sup>2.89</sup> satisfying these conditions.  $\blacktriangledown$

### 2.13.12 Proper nonsimplicial $\mathcal{K}$ , dual, $X$ wide full-rank

Since conically dependent columns can always be removed from  $X$  to construct  $\mathcal{K}$  or to determine  $\mathcal{K}^*$  [413], then assume we are given a set of  $N$  conically independent generators (§2.10) of an arbitrary proper polyhedral cone  $\mathcal{K}$  in  $\mathbb{R}^n$  arranged columnar in  $X \in \mathbb{R}^{n \times N}$  such that  $N > n$  (wide) and  $\text{rank } X = n$ . Having found formula (425) to determine the dual of a simplicial cone, the easiest way to find a vertex-description of proper dual cone  $\mathcal{K}^*$  is to first decompose  $\mathcal{K}$  into simplicial parts  $\mathcal{K}_i$  so that  $\mathcal{K} = \bigcup \mathcal{K}_i$ .<sup>2.90</sup> Each component simplicial cone in  $\mathcal{K}$  corresponds to some subset of  $n$  linearly independent columns from  $X$ . The key idea, here, is how the extreme directions of the simplicial parts must remain extreme directions of  $\mathcal{K}$ . Finding the dual of  $\mathcal{K}$  amounts to finding the dual of each simplicial part:

<sup>2.89</sup> an optimal dual variable, these optimality conditions are equivalent to the KKT conditions [65, §5.5.3].

<sup>2.90</sup>That proposition presupposes, of course, that we know how to perform simplicial decomposition efficiently; also called “triangulation”. [339] [197, §3.1] [198, §3.1] Existence of multiple simplicial parts means expansion of  $x \in \mathcal{K}$ , like (417), can no longer be unique because number  $N$  of extreme directions in  $\mathcal{K}$  exceeds dimension  $n$  of the space.

**2.13.12.0.1 Theorem.** *Dual cone intersection.*

[366, §2.7]

Suppose proper cone  $\mathcal{K} \subset \mathbb{R}^n$  equals the union of  $M$  simplicial cones  $\mathcal{K}_i$  whose extreme directions all coincide with those of  $\mathcal{K}$ . Then proper dual cone  $\mathcal{K}^*$  is the intersection of  $M$  dual simplicial cones  $\mathcal{K}_i^*$ ; *id est*,

$$\mathcal{K} = \bigcup_{i=1}^M \mathcal{K}_i \Rightarrow \mathcal{K}^* = \bigcap_{i=1}^M \mathcal{K}_i^* \quad (475)$$

◊

**Proof.** For  $X_i \in \mathbb{R}^{n \times n}$ , a complete matrix of linearly independent extreme directions (p.114) arranged columnar, corresponding simplicial  $\mathcal{K}_i$  (§2.12.3.1.1) has vertex-description

$$\mathcal{K}_i = \{X_i c \mid c \succeq 0\} \quad (476)$$

Now suppose,

$$\mathcal{K} = \bigcup_{i=1}^M \mathcal{K}_i = \bigcup_{i=1}^M \{X_i c \mid c \succeq 0\} \quad (477)$$

The union of all  $\mathcal{K}_i$  can be equivalently expressed

$$\mathcal{K} = \left\{ [X_1 \ X_2 \ \cdots \ X_M] \begin{bmatrix} a \\ b \\ \vdots \\ c \end{bmatrix} \mid a, b, \dots, c \succeq 0 \right\} \quad (478)$$

Because extreme directions of the simplices  $\mathcal{K}_i$  are extreme directions of  $\mathcal{K}$  by assumption, then

$$\mathcal{K} = \{[X_1 \ X_2 \ \cdots \ X_M] d \mid d \succeq 0\} \quad (479)$$

by the *extremes theorem* (§2.8.1.1.1). Defining  $X \triangleq [X_1 \ X_2 \ \cdots \ X_M]$  (with any redundant [*sic*] columns optionally removed from  $X$ ), then  $\mathcal{K}^*$  can be expressed ((366), Cone Table **S** p.152)

$$\mathcal{K}^* = \{y \mid X^T y \succeq 0\} = \bigcap_{i=1}^M \{y \mid X_i^T y \succeq 0\} = \bigcap_{i=1}^M \mathcal{K}_i^* \quad (480)$$

♦

To find the extreme directions of the dual cone, first we observe: some facets of each simplicial part  $\mathcal{K}_i$  are common to facets of  $\mathcal{K}$  by assumption, and the union of all those common facets comprises the set of all facets of  $\mathcal{K}$  by design. For any particular proper polyhedral cone  $\mathcal{K}$ , the extreme directions of dual cone  $\mathcal{K}^*$  are respectively orthogonal to the facets of  $\mathcal{K}$ . (§2.13.7.1) Then the extreme directions of the dual cone can be found among inward-normals to facets of the component simplicial cones  $\mathcal{K}_i$ ; those normals are extreme directions of the dual simplicial cones  $\mathcal{K}_i^*$ . From the theorem and Cone Table **S** (p.152),

$$\mathcal{K}^* = \bigcap_{i=1}^M \mathcal{K}_i^* = \bigcap_{i=1}^M \{X_i^{\dagger T} c \mid c \succeq 0\} \quad (481)$$

The set of extreme directions  $\{\Gamma_i^*\}$  for proper dual cone  $\mathcal{K}^*$  is therefore constituted by those conically independent generators, from the columns of all the dual simplicial matrices  $\{X_i^{\dagger T}\}$ , that do not violate discrete definition (366) of  $\mathcal{K}^*$ ;

$$\{\Gamma_1^*, \Gamma_2^*, \dots, \Gamma_N^*\} = \text{c.i.} \left\{ X_i^{\dagger T}(:,j), \ i=1 \dots M, \ j=1 \dots n \mid X_i^{\dagger}(j,:) \Gamma_\ell \geq 0, \ \ell=1 \dots N \right\} \quad (482)$$

where c.i. denotes selection of only the conically independent vectors from the argument set, argument  $(:, j)$  denotes the  $j^{\text{th}}$  column while  $(j, :)$  denotes the  $j^{\text{th}}$  row, and  $\{\Gamma_\ell\}$  constitutes the extreme directions of  $\mathcal{K}$ . Figure 53b (p.113) shows a cone and its dual found via this algorithm.

#### 2.13.12.0.2 Example. Dual of $\mathcal{K}$ nonsimplicial in subspace $\text{aff } \mathcal{K}$ .

Given conically independent generators for pointed closed convex cone  $\mathcal{K}$  in  $\mathbb{R}^4$  arranged columnar in

$$X = [\Gamma_1 \ \Gamma_2 \ \Gamma_3 \ \Gamma_4] = \begin{bmatrix} 1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & -1 \end{bmatrix} \quad (483)$$

having  $\dim \text{aff } \mathcal{K} = \text{rank } X = 3$ , (284) then performing the most inefficient simplicial decomposition in  $\text{aff } \mathcal{K}$  we find

$$\begin{aligned} X_1 &= \begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix}, \quad X_2 = \begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & -1 \end{bmatrix} \\ X_3 &= \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -1 & -1 \end{bmatrix}, \quad X_4 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & -1 \end{bmatrix} \end{aligned} \quad (484)$$

The corresponding dual simplicial cones in  $\text{aff } \mathcal{K}$  have generators respectively columnar in

$$\begin{aligned} 4X_1^{\dagger T} &= \begin{bmatrix} 2 & 1 & 1 \\ -2 & 1 & 1 \\ 2 & -3 & 1 \\ -2 & 1 & -3 \end{bmatrix}, \quad 4X_2^{\dagger T} = \begin{bmatrix} 1 & 2 & 1 \\ -3 & 2 & 1 \\ 1 & -2 & 1 \\ 1 & -2 & -3 \end{bmatrix} \\ 4X_3^{\dagger T} &= \begin{bmatrix} 3 & 2 & -1 \\ -1 & 2 & -1 \\ -1 & -2 & 3 \\ -1 & -2 & -1 \end{bmatrix}, \quad 4X_4^{\dagger T} = \begin{bmatrix} 3 & -1 & 2 \\ -1 & 3 & -2 \\ -1 & -1 & 2 \\ -1 & -1 & -2 \end{bmatrix} \end{aligned} \quad (485)$$

Applying algorithm (482) we get

$$[\Gamma_1^* \ \Gamma_2^* \ \Gamma_3^* \ \Gamma_4^*] = \frac{1}{4} \begin{bmatrix} 1 & 2 & 3 & 2 \\ 1 & 2 & -1 & -2 \\ 1 & -2 & -1 & 2 \\ -3 & -2 & -1 & -2 \end{bmatrix} \quad (486)$$

whose rank is 3, and is the known result;<sup>2.91</sup> a conically independent set of generators for that pointed section of the dual cone  $\mathcal{K}^*$  in  $\text{aff } \mathcal{K}$ ; *id est*,  $\mathcal{K}^* \cap \text{aff } \mathcal{K}$ .  $\square$

#### 2.13.12.0.3 Example. Dual of proper polyhedral $\mathcal{K}$ in $\mathbb{R}^4$ .

Given conically independent generators for a full-dimensional pointed closed convex cone  $\mathcal{K}$

$$X = [\Gamma_1 \ \Gamma_2 \ \Gamma_3 \ \Gamma_4 \ \Gamma_5] = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 \\ -1 & 0 & 1 & 0 & 1 \\ 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & -1 & -1 & 0 \end{bmatrix} \quad (487)$$

---

<sup>2.91</sup>These calculations proceed so as to be consistent with [128, §6]; as if the ambient vector space were proper subspace  $\text{aff } \mathcal{K}$  whose dimension is 3. In that ambient space,  $\mathcal{K}$  may be regarded as a proper cone. Yet that author erroneously states dimension of the ordinary dual cone to be 3; it is, in fact, 4.

we count  $5!/((5-4)!4!) = 5$  component simplices.<sup>2.92</sup> Applying algorithm (482), we find the six extreme directions of dual cone  $\mathcal{K}^*$  (with  $\Gamma_2 = \Gamma_5^*$ )

$$X^* = [\Gamma_1^* \ \Gamma_2^* \ \Gamma_3^* \ \Gamma_4^* \ \Gamma_5^* \ \Gamma_6^*] = \begin{bmatrix} 1 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & -1 & 0 & -1 & 1 \\ 1 & -1 & -1 & 1 & 0 & 0 \end{bmatrix} \quad (488)$$

which means, (§2.13.7.1) this proper polyhedral  $\mathcal{K} = \text{cone}(X)$  has six facets generated by its extreme directions:

$$\left\{ \begin{array}{l} \mathcal{G}(\mathcal{F}_1) \\ \mathcal{G}(\mathcal{F}_2) \\ \mathcal{G}(\mathcal{F}_3) \\ \mathcal{G}(\mathcal{F}_4) \\ \mathcal{G}(\mathcal{F}_5) \\ \mathcal{G}(\mathcal{F}_6) \end{array} \right\} = \left\{ \begin{array}{cccccc} \Gamma_1 & \Gamma_2 & \Gamma_3 & & & \\ \Gamma_1 & \Gamma_2 & & & \Gamma_5 & \\ \Gamma_1 & & & \Gamma_4 & \Gamma_5 & \\ \Gamma_1 & & \Gamma_3 & \Gamma_4 & & \\ \Gamma_1 & & \Gamma_3 & \Gamma_4 & \Gamma_5 & \\ \Gamma_2 & \Gamma_3 & & \Gamma_5 & & \end{array} \right\} \quad (489)$$

whereas dual proper polyhedral cone  $\mathcal{K}^*$  has only five (three-dimensional) facets:

$$\left\{ \begin{array}{l} \mathcal{G}(\mathcal{F}_1^*) \\ \mathcal{G}(\mathcal{F}_2^*) \\ \mathcal{G}(\mathcal{F}_3^*) \\ \mathcal{G}(\mathcal{F}_4^*) \\ \mathcal{G}(\mathcal{F}_5^*) \end{array} \right\} = \left\{ \begin{array}{cccccc} \Gamma_1^* & \Gamma_2^* & \Gamma_3^* & \Gamma_4^* & & \\ \Gamma_1^* & \Gamma_2^* & & & \Gamma_6^* & \\ \Gamma_1^* & & \Gamma_4^* & \Gamma_5^* & \Gamma_6^* & \\ \Gamma_1^* & & \Gamma_3^* & \Gamma_4^* & \Gamma_5^* & \\ \Gamma_2^* & \Gamma_3^* & & \Gamma_5^* & \Gamma_6^* & \end{array} \right\} \quad (490)$$

Six two-dimensional cones, having generators  $\mathcal{G}$  respectively  $\{\Gamma_1^* \ \Gamma_3^*\}$   $\{\Gamma_2^* \ \Gamma_4^*\}$   $\{\Gamma_1^* \ \Gamma_5^*\}$   $\{\Gamma_4^* \ \Gamma_6^*\}$   $\{\Gamma_2^* \ \Gamma_5^*\}$   $\{\Gamma_3^* \ \Gamma_6^*\}$ , are relatively interior to dual facets; so cannot be two-dimensional faces of  $\mathcal{K}^*$  (by Definition 2.6.0.0.3).

We can check this result (488) by reversing the process; we find  $6!/((6-4)!4!) - 3 = 12$  component simplices in the dual cone.<sup>2.93</sup> Applying algorithm (482) to those simplices returns a conically independent set of generators for  $\mathcal{K}$  equivalent to (487).  $\square$

#### 2.13.12.0.4 Exercise. Reaching proper polyhedral cone interior.

Name two extreme directions  $\Gamma_i$  of cone  $\mathcal{K}$  from Example 2.13.12.0.3 whose convex hull passes through that cone's interior. Explain why. Are there two such extreme directions of dual cone  $\mathcal{K}^*$ ?  $\blacktriangledown$

### 2.13.13 coordinates in proper nonsimplicial system

A natural question pertains to whether a theory of unique coordinates, like biorthogonal expansion w.r.t pointed closed convex  $\mathcal{K}$ , is extensible to proper cones whose extreme directions number in excess of ambient spatial dimensionality.

#### 2.13.13.0.1 Theorem. Conic coordinates.

With respect to vector  $v$  in some finite-dimensional Euclidean space  $\mathbb{R}^n$ , define a coordinate  $t_v^*$  of point  $x$  in full-dimensional pointed closed convex cone  $\mathcal{K}$

$$t_v^*(x) \triangleq \sup\{t \in \mathbb{R} \mid x - tv \in \mathcal{K}\} \quad (491)$$

Given points  $x$  and  $y$  in cone  $\mathcal{K}$ , if  $t_v^*(x) = t_v^*(y)$  for each and every extreme direction  $v$  of  $\mathcal{K}$  then  $x = y$ .  $\diamond$

<sup>2.92</sup>There are no linearly dependent combinations of three or four extreme directions in the primal cone.

<sup>2.93</sup>Three combinations of four dual extreme directions are linearly dependent; they belong to the dual facets. But there are no linearly dependent combinations of three dual extreme directions.

Conic coordinate definition (491) acquires its heritage from conditions (380) for generator membership to a smallest face. Coordinate  $t_v^*(c)=0$ , for example, corresponds to unbounded  $\mu$  in (380); indicating, extreme direction  $v$  cannot belong to the smallest face of cone  $\mathcal{K}$  that contains  $c$ .

**2.13.13.0.2 Proof.** Vector  $x - t^*v$  must belong to the cone boundary  $\partial\mathcal{K}$  by definition (491). So there must exist a nonzero vector  $\lambda$  that is inward-normal to a hyperplane supporting cone  $\mathcal{K}$  and containing  $x - t^*v$ ; *id est*, by boundary-membership relation for full-dimensional pointed closed convex cones (§2.13.2)

$$x - t^*v \in \partial\mathcal{K} \Leftrightarrow \exists \lambda \neq \mathbf{0} \ni \langle \lambda, x - t^*v \rangle = 0, \quad \lambda \in \mathcal{K}^*, \quad x - t^*v \in \mathcal{K} \quad (333)$$

where

$$\mathcal{K}^* = \{w \in \mathbb{R}^n \mid \langle v, w \rangle \geq 0 \text{ for all } v \in \mathcal{G}(\mathcal{K})\} \quad (372)$$

is the full-dimensional pointed closed convex dual cone. The set  $\mathcal{G}(\mathcal{K})$ , of possibly infinite cardinality  $N$ , comprises generators for cone  $\mathcal{K}$ ; *e.g.*, its extreme directions which constitute a minimal generating set. If  $x - t^*v$  is nonzero, any such vector  $\lambda$  must belong to the dual cone boundary by conjugate boundary-membership relation

$$\lambda \in \partial\mathcal{K}^* \Leftrightarrow \exists x - t^*v \neq \mathbf{0} \ni \langle \lambda, x - t^*v \rangle = 0, \quad x - t^*v \in \mathcal{K}, \quad \lambda \in \mathcal{K}^* \quad (334)$$

where

$$\mathcal{K} = \{z \in \mathbb{R}^n \mid \langle \lambda, z \rangle \geq 0 \text{ for all } \lambda \in \mathcal{G}(\mathcal{K}^*)\} \quad (371)$$

This description of  $\mathcal{K}$  means: cone  $\mathcal{K}$  is an intersection of halfspaces whose inward-normals are generators of the dual cone. Each and every face of cone  $\mathcal{K}$  (except the cone itself) belongs to a hyperplane supporting  $\mathcal{K}$ . Each and every vector  $x - t^*v$  on the cone boundary must therefore be orthogonal to an extreme direction constituting generators  $\mathcal{G}(\mathcal{K}^*)$  of the dual cone.

To the  $i^{\text{th}}$  extreme direction  $v = \Gamma_i \in \mathbb{R}^n$  of cone  $\mathcal{K}$ , ascribe a coordinate  $t_i^*(x) \in \mathbb{R}$  of  $x$  from definition (491). On domain  $\mathcal{K}$ , the mapping

$$t^*(x) = \begin{bmatrix} t_1^*(x) \\ \vdots \\ t_N^*(x) \end{bmatrix} : \mathbb{R}^n \rightarrow \mathbb{R}^N \quad (492)$$

has no nontrivial nullspace. Because  $x - t^*v$  must belong to  $\partial\mathcal{K}$  by definition, the mapping  $t^*(x)$  is equivalent to a convex problem (separable in index  $i$ ) whose objective (by (333)) is tightly bounded below by 0:

$$\begin{aligned} t^*(x) &\equiv \arg \underset{t \in \mathbb{R}^N}{\text{minimize}} \quad \sum_{i=1}^N \Gamma_{j(i)}^{*\top}(x - t_i \Gamma_i) \\ &\text{subject to} \quad x - t_i \Gamma_i \in \mathcal{K}, \quad i = 1 \dots N \end{aligned} \quad (493)$$

where index  $j \in \mathcal{I}$  is dependent on  $i$  and where (by (371))  $\lambda = \Gamma_j^* \in \mathbb{R}^n$  is an extreme direction of dual cone  $\mathcal{K}^*$  that is normal to a hyperplane supporting  $\mathcal{K}$  and containing  $x - t_i^* \Gamma_i$ . Because extreme-direction cardinality  $N$  for cone  $\mathcal{K}$  is not necessarily the same as for dual cone  $\mathcal{K}^*$ , index  $j$  must be judiciously selected from a set  $\mathcal{I}$ .

To prove injectivity when extreme-direction cardinality  $N > n$  exceeds spatial dimension, we need only show mapping  $t^*(x)$  to be invertible; [146, thm.9.2.3] *id est*,  $x$  is recoverable given  $t^*(x)$ :

$$\begin{aligned} x &= \arg \underset{\tilde{x} \in \mathbb{R}^n}{\text{minimize}} \quad \sum_{i=1}^N \Gamma_{j(i)}^{*\top}(\tilde{x} - t_i^* \Gamma_i) \\ &\text{subject to} \quad \tilde{x} - t_i^* \Gamma_i \in \mathcal{K}, \quad i = 1 \dots N \end{aligned} \quad (494)$$

The feasible set of this nonseparable convex problem is an intersection of translated full-dimensional pointed closed convex cones  $\bigcap_i \mathcal{K} + t_i^* \Gamma_i$ . The objective function's linear part describes movement in normal-direction  $-\Gamma_j^*$  for each of  $N$  hyperplanes. The optimal point of hyperplane intersection is the unique solution  $x$  when  $\{\Gamma_j^*\}$  comprises  $n$  linearly independent normals that come from the dual cone and make the objective vanish. Because the dual cone  $\mathcal{K}^*$  is full-dimensional, pointed, closed, and convex by assumption, there exist  $N$  extreme directions  $\{\Gamma_j^*\}$  from  $\mathcal{K}^* \subset \mathbb{R}^n$  that span  $\mathbb{R}^n$ . So we need simply choose  $N$  spanning dual extreme directions that make the optimal objective vanish. Because such dual extreme directions preexist by (333),  $t^*(x)$  is invertible.

Otherwise, in the case  $N \leq n$ ,  $t^*(x)$  holds coordinates for biorthogonal expansion. Reconstruction of  $x$  is therefore unique. ♦

### 2.13.13.1 reconstruction from conic coordinates

The foregoing proof of the *conic coordinates theorem* is not constructive; it establishes existence of dual extreme directions  $\{\Gamma_j^*\}$  that will reconstruct a point  $x$  from its coordinates  $t^*(x)$  via (494), but does not prescribe the index set  $\mathcal{I}$ . There are at least two computational methods for specifying  $\{\Gamma_{j(i)}^*\}$ : one is combinatorial but sure to succeed, the other is a geometric method that searches for a minimum of a nonconvex function. We describe the latter:

Convex problem (P)

$$(P) \quad \begin{array}{ll} \underset{t \in \mathbb{R}}{\text{maximize}} & t \\ \text{subject to} & x - tv \in \mathcal{K} \end{array} \quad \begin{array}{ll} \underset{\lambda \in \mathbb{R}^n}{\text{minimize}} & \lambda^T x \\ \text{subject to} & \lambda^T v = 1 \\ & \lambda \in \mathcal{K}^* \end{array} \quad (D) \quad (495)$$

is equivalent to definition (491) whereas convex problem (D) is its dual;<sup>2.94</sup> meaning, primal and dual optimal objectives are equal  $t^* = \lambda^{*T} x$  assuming Slater's condition (p.229) is satisfied. Under this assumption of strong duality,  $\lambda^{*T}(x - t^*v) = t^*(1 - \lambda^{*T}v) = 0$ ; which implies, the primal problem is equivalent to

$$\begin{array}{ll} \underset{t \in \mathbb{R}}{\text{minimize}} & \lambda^{*T}(x - tv) \\ \text{subject to} & x - tv \in \mathcal{K} \end{array} \quad (p) \quad (496)$$

while the dual problem is equivalent to

$$\begin{array}{ll} \underset{\lambda \in \mathbb{R}^n}{\text{minimize}} & \lambda^T(x - t^*v) \\ \text{subject to} & \lambda^T v = 1 \\ & \lambda \in \mathcal{K}^* \end{array} \quad (d) \quad (497)$$

Instead given coordinates  $t^*(x)$  and a description of cone  $\mathcal{K}$ , we propose inversion by alternating solution of respective primal and dual problems

---

<sup>2.94</sup>Form a Lagrangian associated with primal problem (P):

$$\begin{aligned} \mathfrak{L}(t, \lambda) &= t + \lambda^T(x - tv) = \lambda^T x + t(1 - \lambda^T v), & \lambda \succeq 0 \\ &\sup_{t \in \mathbb{R}} \mathfrak{L}(t, \lambda) = \lambda^T x, & 1 - \lambda^T v = 0 \end{aligned}$$

Dual variable (Lagrange multiplier [280, p.216])  $\lambda$  generally has a nonnegative sense  $\succeq$  for primal maximization with any cone membership constraint, whereas  $\lambda$  would have a nonpositive sense  $\preceq$  were the primal instead a minimization problem with a cone membership constraint.

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \sum_{i=1}^N \Gamma_i^{*\top} (x - t_i^* \Gamma_i) \\ \text{subject to} & x - t_i^* \Gamma_i \in \mathcal{K}, \quad i = 1 \dots N \end{array} \quad (498)$$

$$\begin{array}{ll} \underset{\Gamma_i^* \in \mathbb{R}^n}{\text{minimize}} & \sum_{i=1}^N \Gamma_i^{*\top} (x^* - t_i^* \Gamma_i) \\ \text{subject to} & \Gamma_i^{*\top} \Gamma_i = 1, \quad i = 1 \dots N \\ & \Gamma_i^* \in \mathcal{K}^*, \quad i = 1 \dots N \end{array} \quad (499)$$

where dual extreme directions  $\Gamma_i^*$  are initialized arbitrarily and ultimately ascertained by the alternation. Convex problems (498) and (499) are iterated until convergence which is guaranteed by virtue of a monotonically nonincreasing real sequence of objective values. Convergence can be fast. The mapping  $t^*(x)$  is uniquely inverted when the necessarily nonnegative objective vanishes; *id est*, when  $\Gamma_i^{*\top} (x^* - t_i^* \Gamma_i) = 0 \quad \forall i$ . Here, a zero objective can occur only at the true solution  $x$ . But this global optimality condition cannot be guaranteed by the alternation because the common objective function, when regarded in both primal  $x$  and dual  $\Gamma_i^*$  variables simultaneously, is generally neither quasiconvex or monotonic. (§3.14.0.0.3)

Conversely, a nonzero objective at convergence is a certificate that inversion was not performed properly. A nonzero objective indicates that a global minimum of a multimodal objective function could not be found by this alternation. That is a flaw in this particular iterative algorithm for inversion; not in theory.<sup>2.95</sup> A numerical remedy is to reinitialize the  $\Gamma_i^*$  to different values.

---

<sup>2.95</sup>The Proof 2.13.13.0.2, that suitable dual extreme directions  $\{\Gamma_j^*\}$  always exist, means that a global optimization algorithm would always find the zero objective of alternation (498) (499); hence, the unique inversion  $x$ . But such an algorithm can be combinatorial.

# Chapter 3

## Geometry of convex functions

*The link between convex sets and convex functions is via the epigraph: A function is convex if and only if its epigraph is a convex set.*

—Werner Fenchel

We limit our treatment of *multidimensional functions*<sup>3.1</sup> to finite-dimensional Euclidean space. Then an icon for a one-dimensional (real) *convex function* is bowl-shaped (Figure 82), whereas the *concave* icon is the inverted bowl; respectively characterized by a unique global minimum and maximum whose existence is assumed. Because of this simple relationship, usage of the term *convexity* is often implicitly inclusive of *concavity*. Despite iconic imagery, the reader is reminded that the set of all convex, concave, quasiconvex, and quasiconcave functions contains the *monotonic functions* [229] [242, §2.3.5].

### 3.1 Convex real and vector-valued function

Vector-valued function

$$f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}^M = \begin{bmatrix} f_1(X) \\ \vdots \\ f_M(X) \end{bmatrix} \quad (500)$$

assigns each  $X$  in its domain  $\text{dom } f$  (a subset of ambient vector space  $\mathbb{R}^{p \times k}$ ) to a specific element [289, p.3] of its range (a subset of  $\mathbb{R}^M$ ). Function  $f(X)$  is *linear* in  $X$  on its domain if and only if, for each and every  $Y, Z \in \text{dom } f$  and  $\alpha, \beta \in \mathbb{R}$

$$f(\alpha Y + \beta Z) = \alpha f(Y) + \beta f(Z) \quad (501)$$

A vector-valued function  $f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}^M$  is *convex* in  $X$  if and only if  $\text{dom } f$  is a convex set and, for each and every  $Y, Z \in \text{dom } f$  and  $0 \leq \mu \leq 1$

$$f(\mu Y + (1 - \mu)Z) \underset{\mathbb{R}_+^M}{\preceq} \mu f(Y) + (1 - \mu)f(Z) \quad (502)$$

---

<sup>3.1</sup> vector- or matrix-valued functions including the real functions. Appendix D, with its tables of first- and second-order gradients, is the practical adjunct to this chapter.

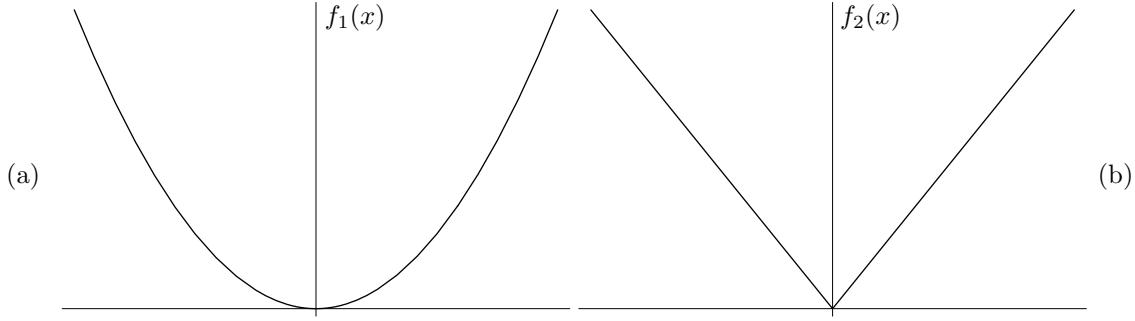


Figure 72: Convex real functions here have a unique minimizer  $x^*$ . For  $x \in \mathbb{R}$ ,  $f_1(x) = x^2 = \|x\|_2^2$  is strictly convex whereas nondifferentiable  $f_2(x) = \sqrt{x^2} = |x| = \|x\|_2$  is convex but not strictly. Strict convexity of a real function is only a sufficient condition for minimizer uniqueness.

As defined, continuity is implied but not differentiability (nor *smoothness*).<sup>3.2</sup> Apparently some, but not all, nonlinear functions are convex. Reversing sense of the inequality flips this definition to concavity. Linear (and affine §3.4)<sup>3.3</sup> functions attain equality in this definition. Linear functions are therefore simultaneously convex and concave.

Vector-valued functions are most often compared (185) as in (502) with respect to the  $M$ -dimensional selfdual nonnegative orthant  $\mathbb{R}_+^M$ , a proper cone.<sup>3.4</sup> In this case, the test prescribed by (502) is simply a comparison on  $\mathbb{R}$  of each *entry*  $f_i$  of a vector-valued function  $f$ . (§2.13.4.2.3) The vector-valued function case is therefore a straightforward generalization of conventional convexity theory for a real function. This conclusion follows from theory of dual generalized inequalities (§2.13.2.0.1) which asserts

$$f \text{ convex w.r.t } \mathbb{R}_+^M \Leftrightarrow w^T f \text{ convex } \forall w \in \mathcal{G}(\mathbb{R}_+^{M^*}) \quad (503)$$

shown by substitution of the defining inequality (502). Discretization allows relaxation (§2.13.4.2.1) of a semiinfinite number of conditions  $\{w \in \mathbb{R}_+^{M^*}\}$  to:

$$\{w \in \mathcal{G}(\mathbb{R}_+^{M^*})\} \equiv \{e_i \in \mathbb{R}^M, i=1 \dots M\} \quad (504)$$

(the standard basis for  $\mathbb{R}^M$  and a minimal set of generators (§2.8.1.2) for  $\mathbb{R}_+^M$ ) from which the stated conclusion follows; *id est*, the test for convexity of a vector-valued function is a comparison on  $\mathbb{R}$  of each entry.

### 3.1.1 strict convexity

When  $f(X)$  instead satisfies, for each and every distinct  $Y$  and  $Z$  in  $\text{dom } f$  and all  $0 < \mu < 1$  on an open interval

$$f(\mu Y + (1 - \mu)Z) \underset{\mathbb{R}_+^M}{\prec} \mu f(Y) + (1 - \mu)f(Z) \quad (505)$$

then we shall say  $f$  is a *strictly convex function*.

<sup>3.2</sup>Figure 72b illustrates a nondifferentiable convex function. Differentiability is certainly not a requirement for optimization of convex functions by numerical methods; *e.g.*, [272].

<sup>3.3</sup>While linear functions are not invariant to translation (offset), convex functions are.

<sup>3.4</sup>Definition of convexity can be broadened to other (not necessarily proper) cones. Referred to in the literature as  $\mathcal{K}$ -convexity, [329]  $\mathbb{R}_+^{M^*}$  (503) generalizes to  $\mathcal{K}^*$ .

Similarly to (503)

$$f \text{ strictly convex w.r.t } \mathbb{R}_+^M \Leftrightarrow w^T f \text{ strictly convex } \forall w \in \mathcal{G}(\mathbb{R}_+^{M^*}) \quad (506)$$

discretization allows relaxation of the semiinfinite number of conditions  $\{w \in \mathbb{R}_+^{M^*}, w \neq \mathbf{0}\}$  (329) to a finite number (504). More tests for strict convexity are given in §3.6.1.0.3, §3.9, and §3.13.0.0.2.

### 3.1.1.1 local/global minimum, uniqueness of solution

A local minimum of any convex real function is also its global minimum. In fact, any convex real function  $f(X)$  has one minimum value over any convex subset of its domain. [338, p.123] Yet solution to some convex optimization problem is, in general, not unique; *id est*, given minimization of a convex real function over some convex feasible set  $\mathcal{C}$

$$\begin{aligned} & \underset{X}{\text{minimize}} && f(X) \\ & \text{subject to} && X \in \mathcal{C} \end{aligned} \quad (507)$$

any *optimal solution*  $X^*$  comes from a convex set of optimal solutions

$$X^* \in \{X \mid f(X) = \inf_{Y \in \mathcal{C}} f(Y)\} \subseteq \mathcal{C} \quad (508)$$

But a strictly convex real function has a unique minimizer  $X^*$ ; *id est*, for the optimal solution set in (508) to be a single point, it is sufficient (Figure 72) that  $f(X)$  be a strictly convex real<sup>3.5</sup> function and set  $\mathcal{C}$  convex. But strict convexity is not necessary for minimizer uniqueness: existence of any strictly supporting hyperplane to a convex function epigraph (p.169, §3.5) at its minimum over  $\mathcal{C}$  is necessary and sufficient for uniqueness.

Quadratic real functions  $x^T A x + b^T x + c$  are convex in  $x$  iff  $A \succeq 0$ . (§3.9.0.0.1) Quadratics characterized by positive definite matrix  $A \succ 0$  are strictly convex and *vice versa*. The vector 2-norm square  $\|x\|^2$  (Euclidean norm square) and Frobenius' norm square  $\|X\|_F^2$ , for example, are strictly convex functions of their respective argument. (Each absolute norm is convex but not strictly.) Figure 72a illustrates a strictly convex real function.

### 3.1.1.2 minimum/minimal element, dual cone characterization

$f(X^*)$  is the *minimum element* of its range if and only if, for each and every  $w \in \text{intr } \mathbb{R}_+^{M^*}$ , it is the unique minimizer of  $w^T f$ . (Figure 73) [65, §2.6.3]

If  $f(X^*)$  is a *minimal element* of its range, then there exists a nonzero  $w \in \mathbb{R}_+^{M^*}$  such that  $f(X^*)$  minimizes  $w^T f$ . If  $f(X^*)$  minimizes  $w^T f$  for some  $w \in \text{intr } \mathbb{R}_+^{M^*}$ , conversely, then  $f(X^*)$  is a minimal element of its range.

#### 3.1.1.2.1 Exercise. Cone of convex functions.

Prove that relation (503) implies: The set of all convex vector-valued functions forms a convex cone in some space. Indeed, any nonnegatively weighted sum of convex functions remains convex. So trivial function  $f = \mathbf{0}$  is convex. Relatively interior to each face of this cone are the strictly convex functions of corresponding dimension.<sup>3.6</sup> How do convex real functions fit into this cone? Where are the affine functions? ▼

<sup>3.5</sup>It is more customary to consider only a real function for the objective of a convex optimization problem because vector- or matrix-valued functions can introduce ambiguity into the optimal objective value. (§2.7.2.2, §3.1.1.2) Study of multidimensional objective functions is called *multicriteria-* [360] or *multiobjective- or vector-optimization*.

<sup>3.6</sup>Strict case excludes cone's point at origin and zero weighting.

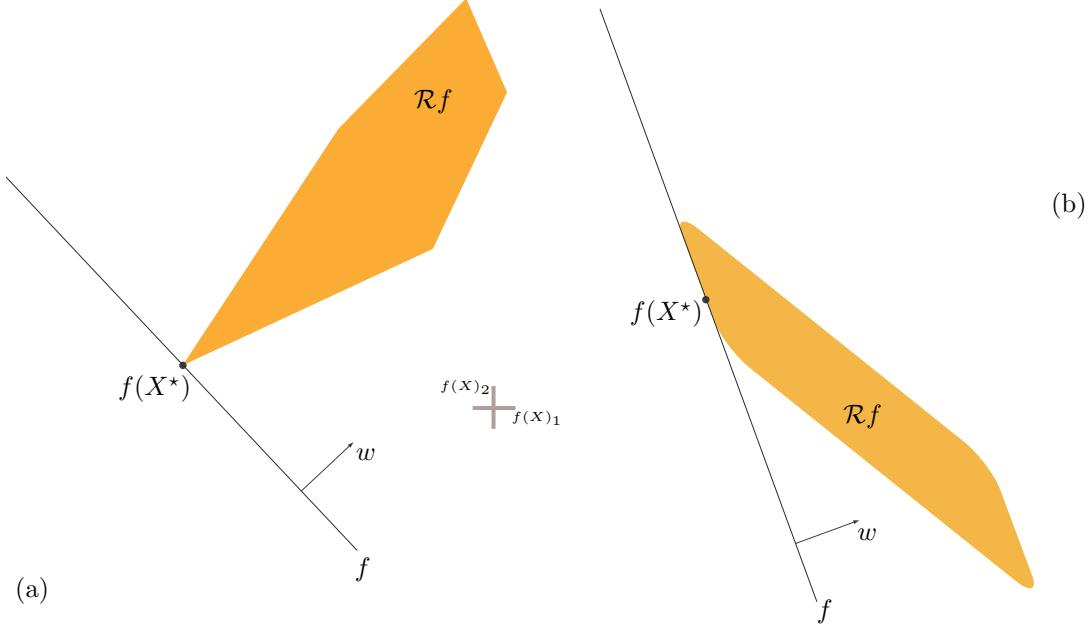


Figure 73: (confer Figure 43) Function range is convex for a convex problem.

(a) Point  $f(X^*)$  is the unique minimum element of function range  $\mathcal{R}f$ .

(b) Point  $f(X^*)$  is a minimal element of depicted range.

(Cartesian axes drawn for reference.)

### 3.1.1.2.2 Example. Conic origins of Lagrangian.

The cone of convex functions, implied by membership relation (503), provides foundation for what is known as a *Lagrangian function*. [282, p.398] [312] Consider a conic optimization problem, for proper cone  $\mathcal{K}$  and affine subset  $\mathcal{A}$

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{subject to} && g(x) \succeq_{\mathcal{K}} 0 \\ & && h(x) \in \mathcal{A} \end{aligned} \tag{509}$$

A Cartesian product of convex functions remains convex, so we may write

$$\begin{bmatrix} f \\ g \\ h \end{bmatrix} \text{ convex w.r.t } \begin{bmatrix} \mathbb{R}_+^M \\ \mathcal{K} \\ \mathcal{A} \end{bmatrix} \Leftrightarrow [w^T \lambda^T \nu^T] \begin{bmatrix} f \\ g \\ h \end{bmatrix} \text{ convex } \forall \begin{bmatrix} w \\ \lambda \\ \nu \end{bmatrix} \in \begin{bmatrix} \mathbb{R}_+^{M^*} \\ \mathcal{K}^* \\ \mathcal{A}^\perp \end{bmatrix} \tag{510}$$

where  $\mathcal{A}^\perp$  is a normal cone to  $\mathcal{A}$ . A normal cone to an affine subset is the orthogonal complement of its parallel subspace (§E.10.3.2.1).

Membership relation (510) holds because of equality for  $h$  in convexity criterion (502) and because normal-cone membership relation (456), given point  $a \in \mathcal{A}$ , becomes

$$h \in \mathcal{A} \Leftrightarrow \langle \nu, h - a \rangle = 0 \text{ for all } \nu \in \mathcal{A}^\perp \tag{511}$$

In other words: all affine functions are convex (with respect to any given proper cone), all convex functions are translation invariant, whereas any affine function must satisfy (511).

A real Lagrangian arises from the scalar term in (510); *id est*,

$$\mathfrak{L} \triangleq [w^T \lambda^T \nu^T] \begin{bmatrix} f \\ g \\ h \end{bmatrix} = w^T f + \lambda^T g + \nu^T h \tag{512}$$

□

## 3.2 Practical norm functions, absolute value

To some mathematicians, “*all norms on  $\mathbb{R}^n$  are equivalent*” [181, p.53]; meaning, ratios of different norms are bounded above and below by finite predetermined constants. But to statisticians and engineers, all norms are certainly not created equal; as evidenced by the *compressed sensing (sparsity)* revolution, begun in 2004, whose focus is predominantly the argument of a 1-norm minimization.

A *norm* on  $\mathbb{R}^n$  is a convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  satisfying: for  $x, y \in \mathbb{R}^n$ ,  $\alpha \in \mathbb{R}$  [254, p.59] [181, p.52]

1.  $f(x) \geq 0$  ( $f(x) = 0 \Leftrightarrow x = 0$ ) (nonnegativity)
2.  $f(x + y) \leq f(x) + f(y)$  3.7 (triangle inequality)
3.  $f(\alpha x) = |\alpha|f(x)$  (nonnegative homogeneity)

Convexity follows by properties 2 and 3. Most useful are 1-, 2-, and  $\infty$ -norm:

$$\begin{aligned} \|x\|_1 &= \underset{t \in \mathbb{R}^n}{\text{minimize}} \quad \mathbf{1}^T t \\ &\text{subject to} \quad -t \preceq x \preceq t \end{aligned} \tag{513}$$

where  $|x| = t^*$  (entrywise absolute value equals optimal  $t$ ).<sup>3.8</sup>

$$\begin{aligned} \|x\|_1 &= \underset{\alpha \in \mathbb{R}^n, \beta \in \mathbb{R}^n}{\text{minimize}} \quad \mathbf{1}^T (\alpha + \beta) \\ &\text{subject to} \quad \alpha, \beta \succeq 0 \\ &\quad x = \alpha - \beta \end{aligned} \tag{514}$$

where  $|x| = \alpha^* + \beta^*$  because of complementarity  $\alpha^* \mathbf{1}^T \beta^* = 0$  at optimality.

$$\begin{aligned} \|x\|_2 &= \underset{t \in \mathbb{R}}{\text{minimize}} \quad t \\ &\text{subject to} \quad \begin{bmatrix} tI & x \\ x^T & t \end{bmatrix} \succeq_{\mathbb{S}_+^{n+1}} 0 \end{aligned} \tag{515}$$

where  $\|x\|_2 = \|x\| \triangleq \sqrt{x^T x} = t^*$ .

$$\begin{aligned} \|x\|_\infty &= \underset{t \in \mathbb{R}}{\text{minimize}} \quad t \\ &\text{subject to} \quad -t \mathbf{1} \preceq x \preceq t \mathbf{1} \end{aligned} \tag{516}$$

where  $\max\{|x_i|, i=1 \dots n\} = t^*$  because  $\|x\|_\infty = \max\{|x_i|\} \leq t \Leftrightarrow |x| \preceq t \mathbf{1}$ ; absolute value  $|x|$  inequality, in this sense, describing a norm ball. (513) (514) and (516) represent linear programs, (515) is a semidefinite program.

Each of these norms is a monotonic real function on a nonnegative orthant. But over any arbitrary convex set  $\mathcal{C}$ , given vector constant  $y$  or matrix constant  $Y$

$$\arg \inf_{x \in \mathcal{C}} \|x - y\| = \arg \inf_{x \in \mathcal{C}} \|x - y\|^2 \tag{517}$$

$$\arg \inf_{X \in \mathcal{C}} \|X - Y\| = \arg \inf_{X \in \mathcal{C}} \|X - Y\|^2 \tag{518}$$

<sup>3.7</sup>  $\|x + y\| \leq \|x\| + \|y\|$  for any norm, with equality iff  $x = \kappa y$  where  $\kappa \geq 0$ .

<sup>3.8</sup> Vector  $\mathbf{1}$  may be replaced with any positive [*sic*] vector to get absolute value, theoretically, although  $\mathbf{1}$  provides the 1-norm.

are unconstrained convex problems for any convex norm and any affine transformation of variable. (But equality would not hold for, instead, a sum of norms; *e.g.* §5.4.2.2.4.) Optimal solution is norm dependent: [65, p.297]

$$\begin{array}{lll} \underset{\substack{x \in \mathbb{R}^n \\ \text{subject to} \\ x \in \mathcal{C}}}{\text{minimize}} & \|x\|_1 & \underset{\substack{x \in \mathbb{R}^n, t \in \mathbb{R} \\ \text{subject to} \\ -t \leq x \leq t \\ x \in \mathcal{C}}}{\text{minimize}} \quad \mathbf{1}^T t \\ & & \end{array} \quad (519)$$

$$\begin{array}{lll} \underset{\substack{x \in \mathbb{R}^n \\ \text{subject to} \\ x \in \mathcal{C}}}{\text{minimize}} & \|x\|_2 & \underset{\substack{x \in \mathbb{R}^n, t \in \mathbb{R} \\ \text{subject to} \\ \left[ \begin{array}{cc} tI & x \\ x^T & t \end{array} \right] \succeq 0 \\ x \in \mathcal{C}}}{\text{minimize}} \quad t \\ & & \end{array} \quad (520)$$

$$\begin{array}{lll} \underset{\substack{x \in \mathbb{R}^n \\ \text{subject to} \\ x \in \mathcal{C}}}{\text{minimize}} & \|x\|_\infty & \underset{\substack{x \in \mathbb{R}^n, t \in \mathbb{R} \\ \text{subject to} \\ -t\mathbf{1} \leq x \leq t\mathbf{1} \\ x \in \mathcal{C}}}{\text{minimize}} \quad t \\ & & \end{array} \quad (521)$$

In  $\mathbb{R}^n$ :  $\|x\|_1$  represents length measured along a grid (*taxicab distance*, Figure 84),  $\|x\|_2$  is Euclidean length,  $\|x\|_\infty$  is maximum |coordinate| (*peak magnitude*).

$$\begin{array}{lll} \underset{\substack{x \in \mathbb{R}^n \\ \text{subject to} \\ x \in \mathcal{C}}}{\text{minimize}} & \|x\|_1 & \underset{\substack{\alpha \in \mathbb{R}^n, \beta \in \mathbb{R}^n \\ \text{subject to} \\ \alpha, \beta \succeq 0 \\ x = \alpha - \beta \\ x \in \mathcal{C}}}{\text{minimize}} \quad \mathbf{1}^T(\alpha + \beta) \\ & & \end{array} \quad (522)$$

These foregoing problems (513)-(522) are convex whenever set  $\mathcal{C}$  is. Their equivalence transformations make objectives smooth.

### 3.2.0.0.1 Example. Projecting the origin, on an affine subset, in 1-norm.

In (2064) we interpret *least norm* solution to linear system  $Ax = b$  as orthogonal projection of the origin  $\mathbf{0}$  on affine subset  $\mathcal{A} = \{x \in \mathbb{R}^n \mid Ax = b\}$  where  $A \in \mathbb{R}^{m \times n}$  is wide full-rank. Suppose, instead of the Euclidean metric, we use taxicab distance to do projection. Then the least 1-norm problem is stated, for  $b \in \mathcal{R}(A)$

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|x\|_1 \\ \text{subject to} & Ax = b \end{array} \quad (523)$$

**a.k.a compressed sensing problem.** Optimal solution can be interpreted as an *oblique projection* of the origin on  $\mathcal{A}$  simply because the Euclidean metric is not employed. This problem statement sometimes returns optimal  $x^*$  having minimal cardinality; which can be explained intuitively with reference to Figure 74: [20]

Projection of the origin, in 1-norm, on affine subset  $\mathcal{A}$  is equivalent to maximization (in this case) of the 1-norm ball  $\mathcal{B}_1$  until it kisses  $\mathcal{A}$ ; rather, a kissing point in  $\mathcal{A}$  achieves the distance in 1-norm from the origin to  $\mathcal{A}$ . For the example illustrated ( $m=1$ ,  $n=3$ ), it appears that a vertex of the ball will be first to touch  $\mathcal{A}$ . 1-norm ball vertices in  $\mathbb{R}^3$  represent nontrivial points of minimal cardinality 1, whereas edges represent cardinality 2, while relative interiors of facets represent maximal cardinality 3. By reorienting affine subset  $\mathcal{A}$  so it were parallel to an edge or facet, it becomes evident as we expand or contract the ball that a kissing point is not necessarily unique.<sup>3.9</sup>

---

<sup>3.9</sup>This is unlike the case for the Euclidean ball (2064) where minimum-distance projection on a convex set is unique (§E.9); all kissable faces of the Euclidean ball are single points (vertices).

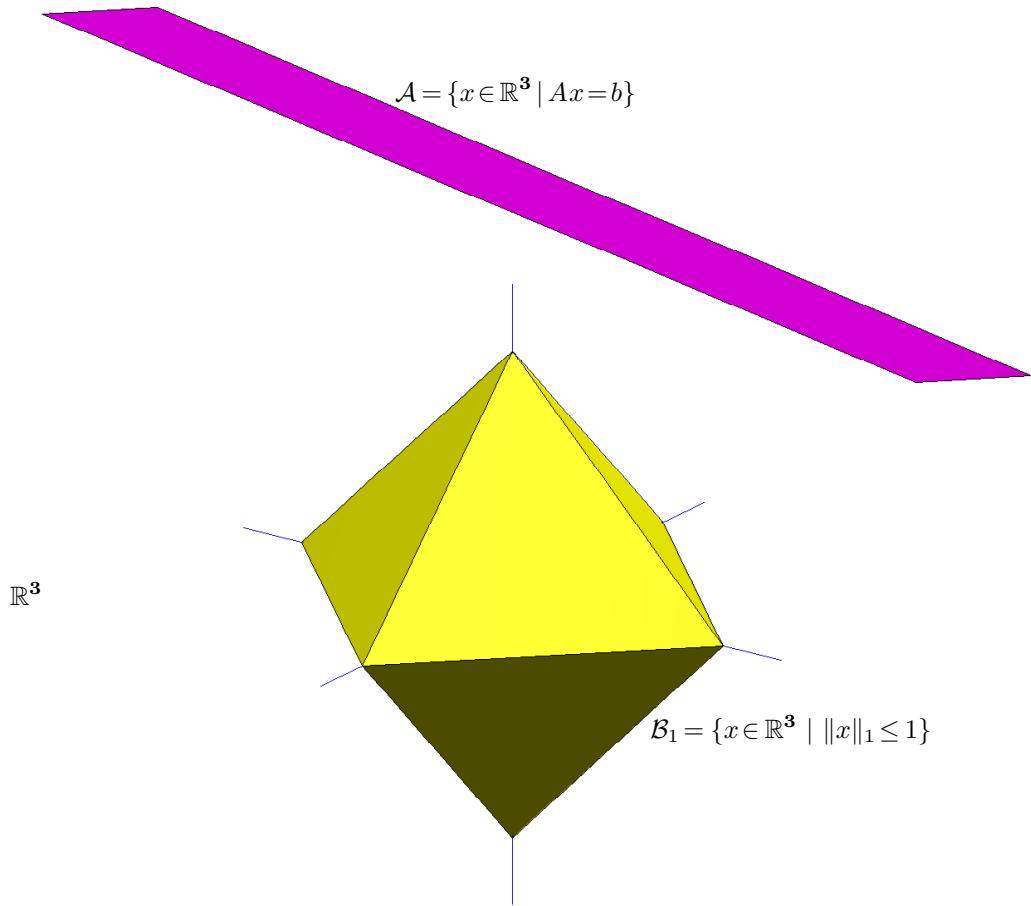
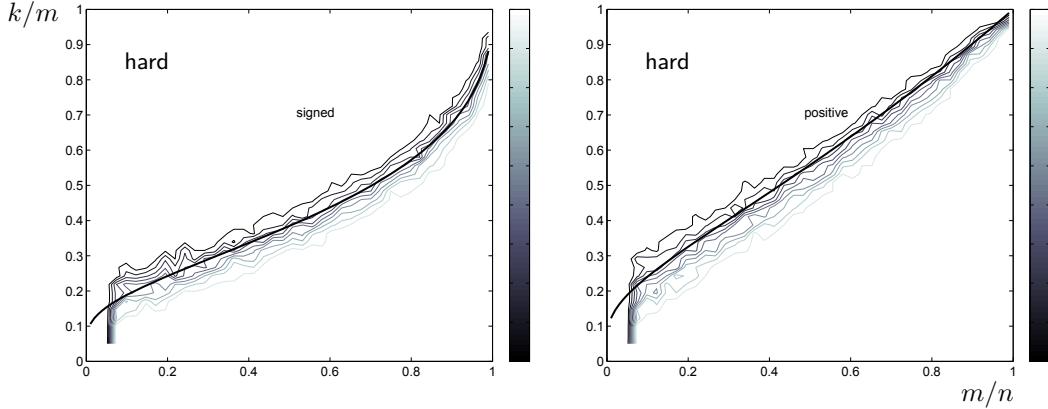


Figure 74: 1-norm ball  $\mathcal{B}_1$  is convex hull of all cardinality-1 vectors of unit norm (its vertices). Ball boundary contains all points equidistant from origin in 1-norm. (Cartesian axes drawn for reference.) Plane  $\mathcal{A}$  is overhead (drawn truncated). If 1-norm ball is expanded until it *kisses*  $\mathcal{A}$  (intersects ball only at boundary), then distance (in 1-norm) from origin to  $\mathcal{A}$  is achieved. Euclidean ball would be spherical in this dimension. Only were  $\mathcal{A}$  parallel to two axes could there be a minimal cardinality least Euclidean norm solution. Yet 1-norm ball offers infinitely many, but not all,  $\mathcal{A}$ -orientations resulting in a minimal cardinality solution. (1-norm ball is an octahedron in this dimension while  $\infty$ -norm ball is a cube.)



$$(523) \quad \begin{aligned} & \underset{x}{\text{minimize}} && \|x\|_1 \\ & \text{subject to} && Ax = b \end{aligned}$$

$$(528) \quad \begin{aligned} & \underset{x}{\text{minimize}} && \|x\|_1 \\ & \text{subject to} && Ax = b \\ & && x \succeq 0 \end{aligned}$$

Figure 75: (confer Figure 113) Exact recovery transition: Respectively signed [135] [137] or positive [142] [140] [141] solutions  $x$ , to  $Ax=b$  with sparsity  $k$  below thick curve, are recoverable. For Gaussian random matrix  $A \in \mathbb{R}^{m \times n}$ , thick curve demarcates *phase transition* in ability to find sparsest solution  $x$  by linear programming. These results empirically reproduced in [39].

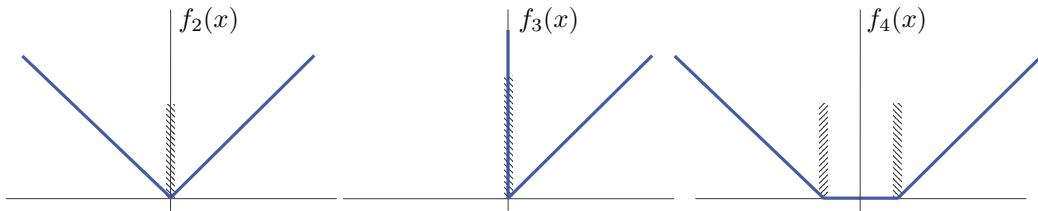


Figure 76: Under 1-norm  $f_2(x)$ , histogram (hatched) of residual amplitudes  $Ax-b$  exhibits predominant accumulation of zero-residuals. Nonnegatively constrained 1-norm  $f_3(x)$  from (528) accumulates more zero-residuals than  $f_2(x)$ . Under norm  $f_4(x)$  (not discussed), histogram would exhibit predominant accumulation of (nonzero) residuals at gradient discontinuities.

The 1-norm ball in  $\mathbb{R}^n$  has  $2^n$  facets and  $2n$  vertices.<sup>3.10</sup> For  $n > 0$

$$\mathcal{B}_1 = \{x \in \mathbb{R}^n \mid \|x\|_1 \leq 1\} = \text{conv}\{\|x \in \mathbb{R}^n\| = 1 \mid \text{card } x = 1\} = \text{conv}\{\pm e_i \in \mathbb{R}^n, i = 1 \dots n\} \quad (524)$$

is a vertex-description of the unit 1-norm ball. Maximization of the 1-norm ball, until it kisses  $\mathcal{A}$ , is equivalent to minimization of the 1-norm ball until it no longer intersects  $\mathcal{A}$ . Then projection of the origin on affine subset  $\mathcal{A}$  is

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \|x\|_1 \\ \text{subject to} & Ax = b \end{array} \equiv \begin{array}{ll} \underset{c \in \mathbb{R}, x \in \mathbb{R}^n}{\text{minimize}} & c \\ \text{subject to} & x \in c\mathcal{B}_1 \\ & Ax = b \end{array} \quad (525)$$

where

$$c\mathcal{B}_1 = \{[I \in \mathbb{R}^{n \times n} \ -I \in \mathbb{R}^{n \times n}]a \mid a^T \mathbf{1} = c, a \succeq 0\} \quad (526)$$

which is the convex hull of 1-norm ball vertices. Then (525) is equivalent to

$$\begin{array}{ll} \underset{c \in \mathbb{R}, x \in \mathbb{R}^n, a \in \mathbb{R}^{2n}}{\text{minimize}} & c \\ \text{subject to} & x = [I \ -I]a \\ & a^T \mathbf{1} = c \\ & a \succeq 0 \\ & Ax = b \end{array} \equiv \begin{array}{ll} \underset{a \in \mathbb{R}^{2n}}{\text{minimize}} & \|a\|_1 \\ \text{subject to} & [A \ -A]a = b \\ & a \succeq 0 \end{array} \quad (527)$$

where  $x^* = [I \ -I]a^*$ . (confer (522)) Significance of this result:

- (confer p.330) Any vector 1-norm minimization problem may have its variable replaced with a nonnegative variable of the same optimal cardinality but twice the length.

All other things being equal, nonnegative variables are easier to solve for sparse solutions. (Figure 75, Figure 76, Figure 113) The compressed sensing problem (523) becomes easier to interpret; e.g., for  $A \in \mathbb{R}^{m \times n}$

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|x\|_1 \\ \text{subject to} & Ax = b \\ & x \succeq 0 \end{array} \equiv \begin{array}{ll} \underset{x}{\text{minimize}} & \mathbf{1}^T x \\ \text{subject to} & Ax = b \\ & x \succeq 0 \end{array} \quad (528)$$

movement of a hyperplane (Figure 29, Figure 33) over a bounded polyhedron always has a *vertex solution* [103, p.158]. Or vector  $b$  might lie on the relative boundary of a pointed polyhedral cone  $\mathcal{K} = \{Ax \mid x \succeq 0\}$ . In the latter case, we find practical application of the smallest face  $\mathcal{F}$  containing  $b$  from §2.13.5 to remove all columns of matrix  $A$  not belonging to  $\mathcal{F}$ ; because those columns correspond to 0-entries in vector  $x$  (and *vice versa*).  $\square$

### 3.2.0.0.2 Exercise. Combinatorial optimization.

A device commonly employed to relax combinatorial problems is to arrange desirable solutions at vertices of bounded polyhedra; e.g., the permutation matrices of dimension  $n$ , which are factorial in number, are the extreme points of a polyhedron in the nonnegative orthant described by an intersection of  $2n$  hyperplanes (§2.3.2.0.4). Minimizing a linear objective function over a bounded polyhedron is a convex problem (a linear program) that always has an optimal solution residing at a vertex.

What about minimizing other functions? Given some *nonsingular* matrix  $A$ , describe three circumstances geometrically under which there are likely to exist vertex solutions to

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \|Ax\|_1 \\ \text{subject to} & x \in \mathcal{P} \end{array} \quad (529)$$

---

<sup>3.10</sup>The  $\infty$ -norm ball in  $\mathbb{R}^n$  has  $2n$  facets and  $2^n$  vertices.

optimized over some bounded polyhedron  $\mathcal{P}$ .<sup>3.11</sup> ▼

### 3.2.1 $k$ smallest entries

Sum of the  $k$  smallest entries of  $x \in \mathbb{R}^n$  is the optimal objective value from: for  $1 \leq k \leq n$

$$\begin{array}{lll} \sum_{i=n-k+1}^n \pi(x)_i = \underset{y \in \mathbb{R}^n}{\text{minimize}} & x^T y \\ \text{subject to} & 0 \preceq y \preceq \mathbf{1} \\ & \mathbf{1}^T y = k \end{array} \equiv \begin{array}{lll} \sum_{i=n-k+1}^n \pi(x)_i = \underset{z \in \mathbb{R}^n, t \in \mathbb{R}}{\text{maximize}} & k t + \mathbf{1}^T z \\ \text{subject to} & x \succeq t \mathbf{1} + z \\ & z \preceq 0 \end{array} \quad (530)$$

which are dual linear programs, where  $\pi(x)_1 = \max\{x_i, i=1 \dots n\}$  where  $\pi$  is a nonlinear permutation-operator sorting its vector argument into nonincreasing order. Finding  $k$  smallest entries of an  $n$ -length vector  $x$  is expressible as an infimum of  $n!/(k!(n-k)!)$  linear functions of  $x$ . The sum  $\sum \pi(x)_i$  is therefore a concave function of  $x$ ; in fact, monotonic (§3.6.1.0.1) in this instance.

### 3.2.2 $k$ largest entries

Sum of the  $k$  largest entries of  $x \in \mathbb{R}^n$  is the optimal objective value from: [65, exer.5.19]

$$\begin{array}{lll} \sum_{i=1}^k \pi(x)_i = \underset{y \in \mathbb{R}^n}{\text{maximize}} & x^T y \\ \text{subject to} & 0 \preceq y \preceq \mathbf{1} \\ & \mathbf{1}^T y = k \end{array} \equiv \begin{array}{lll} \sum_{i=1}^k \pi(x)_i = \underset{z \in \mathbb{R}^n, t \in \mathbb{R}}{\text{minimize}} & k t + \mathbf{1}^T z \\ \text{subject to} & x \preceq t \mathbf{1} + z \\ & z \succeq 0 \end{array} \quad (531)$$

which are dual linear programs. Finding  $k$  largest entries of an  $n$ -length vector  $x$  is expressible as a supremum of  $n!/(k!(n-k)!)$  linear functions of  $x$ . (Figure 78) The summation is therefore a convex function (and monotonic in this instance, §3.6.1.0.1).

#### 3.2.2.1 $k$ -largest norm

Let  $\Pi x$  be a permutation of entries  $x_i$  such that their absolute value becomes arranged in nonincreasing order:  $|\Pi x|_1 \geq |\Pi x|_2 \geq \dots \geq |\Pi x|_n$ . Sum of the  $k$  largest entries of  $|x| \in \mathbb{R}^n$  is a norm, by properties of vector norm (§3.2), and is the optimal objective value of a linear program:

$$\begin{aligned} \|x\|_k^n &\triangleq \sum_{i=1}^k |\Pi x|_i = \sum_{i=1}^k \pi(|x|)_i = \underset{z \in \mathbb{R}^n, t \in \mathbb{R}}{\text{minimize}} & k t + \mathbf{1}^T z \\ &\text{subject to} & -t \mathbf{1} - z \preceq x \preceq t \mathbf{1} + z \\ && z \succeq 0 \\ &= \sup_{i \in \mathcal{I}} \left\{ a_i^T x \mid \begin{array}{l} a_{ij} \in \{-1, 0, 1\} \\ \text{card } a_i = k \end{array} \right\} = \underset{y_1, y_2 \in \mathbb{R}^n}{\text{maximize}} & (y_1 - y_2)^T x \\ &\text{subject to} & 0 \preceq y_1 \preceq \mathbf{1} \\ && 0 \preceq y_2 \preceq \mathbf{1} \\ && (y_1 + y_2)^T \mathbf{1} = k \end{aligned} \quad (532)$$

where the norm subscript derives from a binomial coefficient  $\binom{n}{k}$ , and

---

<sup>3.11</sup>Hint: Suppose, for example,  $\mathcal{P}$  belongs to an orthant and  $A$  were orthogonal. Begin with  $A = I$  and apply level sets of the objective, as in Figure 71 and Figure 74. Or rewrite the problem as a linear program like (519) and (521) but in a composite variable  $\begin{bmatrix} x \\ t \end{bmatrix} \leftarrow y$ .

$$\begin{aligned}\|x\|_n &= \|x\|_1 \\ \|x\|_1 &= \|x\|_\infty \\ \|x\|_k &= \|\pi(|x|)_{1:k}\|_1\end{aligned}\tag{533}$$

Sum of  $k$  largest absolute entries of an  $n$ -length vector  $x$  is expressible as a supremum of  $2^k n!/(k!(n-k)!)$  linear functions of  $x$ ; (Figure 78) hence, this norm is convex in  $x$ . [65, exer.6.3e]

$$\begin{array}{lll} \text{minimize}_{x \in \mathbb{R}^n} & \|x\|_k^n & \equiv \begin{array}{ll} \text{minimize}_{z \in \mathbb{R}^n, t \in \mathbb{R}, x \in \mathbb{R}^n} & kt + \mathbf{1}^T z \\ \text{subject to} & -t\mathbf{1} - z \preceq x \preceq t\mathbf{1} + z \\ & z \succeq 0 \\ & x \in \mathcal{C} \end{array} \\ \text{subject to} & x \in \mathcal{C} & \end{array} \tag{534}$$

### 3.2.2.1.1 Exercise. Polyhedral epigraph of $k$ -largest norm.

Make those  $\text{card } \mathcal{I} = 2^k n!/(k!(n-k)!)$  linear functions explicit for  $\|x\|_2$  and  $\|x\|_1$  on  $\mathbb{R}^2$  and  $\|x\|_3$  on  $\mathbb{R}^3$ . Plot  $\|x\|_2$  and  $\|x\|_1$  in three dimensions. ▼

### 3.2.2.1.2 Exercise. Norm strict convexity.

Which of the vector norms  $\|x\|_k$ ,  $\|x\|_1$ ,  $\|x\|_2$ ,  $\|x\|_\infty$  become strictly convex when squared? Do they remain strictly convex when raised to third and higher whole powers? ▼

### 3.2.2.1.3 Example. Compressed sensing problem.

Conventionally posed as convex problem (523), we showed: the compressed sensing problem can always be posed equivalently with a nonnegative variable as in convex statement (528). The 1-norm predominantly appears in the literature because it is convex, it inherently minimizes cardinality under some technical conditions, [75] and because the desirable 0-norm is intractable.

Assuming a cardinality- $k$  solution exists, the compressed sensing problem may be written as a difference of two convex functions: for  $A \in \mathbb{R}^{m \times n}$

$$\begin{array}{lll} \text{minimize}_{x \in \mathbb{R}^n} & \|x\|_1 - \|x\|_k^n & \text{find } x \in \mathbb{R}^n \\ \text{subject to} & Ax = b & \equiv \text{subject to } \begin{array}{l} Ax = b \\ x \succeq 0 \\ \|x\|_0 \leq k \end{array} \end{array} \tag{535}$$

which is a nonconvex statement, a minimization of  $n-k$  smallest entries of variable vector  $x$ , minimization of a concave function on  $\mathbb{R}_+^n$  (§3.2.1) [343, §32]; but a statement of compressed sensing more precise than (528) because of its equivalence to 0-norm.  $\|x\|_k^n$  is the convex  $k$ -largest norm of  $x$  (monotonic on  $\mathbb{R}_+^n$ ) while  $\|x\|_0$  expresses its cardinality (quasiconcave on  $\mathbb{R}_+^n$ ). Global optimality occurs at a zero objective of minimization; *id est*, when the smallest  $n-k$  entries of variable vector  $x$  are zeroed. Under nonnegativity constraint, this compressed sensing problem (535) becomes the same as

$$\begin{array}{lll} \text{minimize}_{z(x), x \in \mathbb{R}^n} & (\mathbf{1} - z)^T x & \\ \text{subject to} & Ax = b & \\ & x \succeq 0 & \end{array} \tag{536}$$

where

$$\left. \begin{array}{l} \mathbf{1} = \nabla \|x\|_1 = \nabla \mathbf{1}^T x \\ z = \nabla \|x\|_k^n = \nabla z^T x \end{array} \right\}, \quad x \succeq 0 \tag{537}$$

where gradient of  $k$ -largest norm is an optimal solution to a convex problem:

$$\left. \begin{array}{l} \|x\|_k = \max_{y \in \mathbb{R}^n} y^T x \\ \text{subject to } 0 \preceq y \preceq \mathbf{1} \\ y^T \mathbf{1} = k \\ \nabla \|x\|_k = \arg \max_{y \in \mathbb{R}^n} y^T x \\ \text{subject to } 0 \preceq y \preceq \mathbf{1} \\ y^T \mathbf{1} = k \end{array} \right\}, \quad x \succeq 0 \quad (538)$$

□

### 3.2.2.1.4 Exercise. $k$ -largest norm gradient.

Prove (537). Is  $\nabla \|x\|_k$  unique? Find  $\nabla \|x\|_1$  and  $\nabla \|x\|_k$  on  $\mathbb{R}^n$ . [3.12](#)

▼

## 3.2.3 clipping

Zeroing negative vector entries under 1-norm is accomplished:

$$\begin{aligned} \|x_+\|_1 &= \min_{t \in \mathbb{R}^n} \mathbf{1}^T t \\ \text{subject to } &x \preceq t \\ &0 \preceq t \end{aligned} \quad (539)$$

where, for  $x = [x_i, i=1 \dots n] \in \mathbb{R}^n$

$$x_+ \triangleq t^* = \begin{bmatrix} x_i, & x_i \geq 0 \\ 0, & x_i < 0 \end{bmatrix}, \quad i=1 \dots n = \frac{1}{2}(x + |x|) \quad (540)$$

(clipping)

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \|x_+\|_1 &\equiv \min_{\substack{x \in \mathbb{R}^n, t \in \mathbb{R}^n \\ \text{subject to } x \in \mathcal{C}}} \mathbf{1}^T t \\ &\text{subject to } x \preceq t \\ &0 \preceq t \\ &x \in \mathcal{C} \end{aligned} \quad (541)$$

## 3.3 Powers, roots, and inverted functions

A given function  $f$  is convex iff  $-f$  is concave. Both functions are loosely referred to as *convex* since  $-f$  is simply  $f$  inverted about the abscissa axis, and minimization of  $f$  is equivalent to maximization of  $-f$ .

A given positive function  $f$  is convex iff  $1/f$  is concave;  $f$  inverted about ordinate 1 is concave. Minimization of  $f$  is maximization of  $1/f$ .

We wish to implement objectives of the form  $x^{-1}$ . Suppose we have a  $2 \times 2$  matrix

$$T \triangleq \begin{bmatrix} x & z \\ z & y \end{bmatrix} \in \mathbb{R}^2 \quad (542)$$

which is positive semidefinite by (1666) when

$$T \succeq 0 \Leftrightarrow x > 0 \text{ and } xy \geq z^2 \quad (543)$$

---

[3.12](#) Hint: §D.2.1.

A polynomial constraint such as this is therefore called a *conic constraint*.<sup>3.13</sup> This means we may formulate convex problems, having inverted variables, as semidefinite programs in Schur-form (§A.4); *e.g.*,

$$\begin{array}{ll} \underset{x \in \mathbb{R}}{\text{minimize}} & x^{-1} \\ \text{subject to} & x > 0 \\ & x \in \mathcal{C} \end{array} \equiv \begin{array}{ll} \underset{x, y \in \mathbb{R}}{\text{minimize}} & y \\ \text{subject to} & \begin{bmatrix} x & 1 \\ 1 & y \end{bmatrix} \succeq 0 \\ & x \in \mathcal{C} \end{array} \quad (544)$$

rather

$$x > 0, \quad y \geq \frac{1}{x} \Leftrightarrow \begin{bmatrix} x & 1 \\ 1 & y \end{bmatrix} \succeq 0 \quad (545)$$

(inverted) For vector  $x = [x_i, i=1 \dots n] \in \mathbb{R}^n$

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \sum_{i=1}^n x_i^{-1} \\ \text{subject to} & x \succ 0 \\ & x \in \mathcal{C} \end{array} \equiv \begin{array}{ll} \underset{x \in \mathbb{R}^n, y \in \mathbb{R}}{\text{minimize}} & y \\ \text{subject to} & \begin{bmatrix} x_i & \sqrt{n} \\ \sqrt{n} & y \end{bmatrix} \succeq 0, \quad i=1 \dots n \\ & x \in \mathcal{C} \end{array} \quad (546)$$

rather

$$x \succ 0, \quad y \geq \text{tr}(\delta(x)^{-1}) \Leftrightarrow \begin{bmatrix} x_i & \sqrt{n} \\ \sqrt{n} & y \end{bmatrix} \succeq 0, \quad i=1 \dots n \quad (547)$$

### 3.3.1 rational exponent

Galtier [169] shows how to implement an objective of the form  $x^\alpha$  for positive  $\alpha$ . He suggests quantizing  $\alpha$  and working instead with that approximation. Choose nonnegative integer  $q$  for adequate quantization of  $\alpha$  like so:

$$\alpha \triangleq \frac{k}{2^q}, \quad k \in \{0, 1, 2 \dots 2^q - 1\} \quad (548)$$

Any  $k$  from that set may be written  $k = \sum_{i=1}^q b_i 2^{i-1}$  where  $b_i \in \{0, 1\}$ . Define vector  $y = [y_i, i=0 \dots q] \in \mathbb{R}^{q+1}$  with  $y_0 = 1$ :

#### 3.3.1.1 positive

Then we have the equivalent semidefinite program for maximizing a concave function  $x^\alpha$ , for quantized  $0 \leq \alpha < 1$

$$\begin{array}{ll} \underset{x \in \mathbb{R}}{\text{maximize}} & x^\alpha \\ \text{subject to} & x > 0 \\ & x \in \mathcal{C} \end{array} \equiv \begin{array}{ll} \underset{x \in \mathbb{R}, y \in \mathbb{R}^{q+1}}{\text{maximize}} & y_q \\ \text{subject to} & \begin{bmatrix} y_{i-1} & y_i \\ y_i & x^{b_i} \end{bmatrix} \succeq 0, \quad i=1 \dots q \\ & x \in \mathcal{C} \end{array} \quad (549)$$

where nonnegativity of  $y_q$  is enforced by maximization; *id est*,

$$x > 0, \quad y_q \leq x^\alpha \Leftrightarrow \begin{bmatrix} y_{i-1} & y_i \\ y_i & x^{b_i} \end{bmatrix} \succeq 0, \quad i=1 \dots q \quad (550)$$

---

<sup>3.13</sup>In this dimension, the convex cone formed from the set of all values  $\{x, y, z\}$  (satisfying constraint (543)) is called a positive semidefinite cone or a *rotated* quadratic, circular, or second-order cone.

**3.3.1.1.1 Example.** *Square root.*

$\alpha = \frac{1}{2}$ . Choose  $q=1$  and  $k=1=2^0$ .

$$\begin{array}{ll} \text{maximize}_{x \in \mathbb{R}} & \sqrt{x} \\ \text{subject to} & x > 0 \\ & x \in \mathcal{C} \end{array} \equiv \begin{array}{ll} \text{maximize}_{x \in \mathbb{R}, y \in \mathbb{R}^2} & y_1 \\ \text{subject to} & \begin{bmatrix} y_0 = 1 & y_1 \\ y_1 & x \end{bmatrix} \succeq 0 \\ & x \in \mathcal{C} \end{array} \quad (551)$$

where

$$x > 0, \quad y_1 \leq \sqrt{x} \Leftrightarrow \begin{bmatrix} 1 & y_1 \\ y_1 & x \end{bmatrix} \succeq 0 \quad (552)$$

□

**3.3.1.2 negative**

It is also desirable to implement an objective of the form  $x^{-\alpha}$  for positive  $\alpha$ . The technique is nearly the same as before: for quantized  $0 \leq \alpha < 1$

$$\begin{array}{ll} \text{minimize}_{x \in \mathbb{R}} & x^{-\alpha} \\ \text{subject to} & x > 0 \\ & x \in \mathcal{C} \end{array} \equiv \begin{array}{ll} \text{minimize}_{x, z \in \mathbb{R}, y \in \mathbb{R}^{q+1}} & z \\ \text{subject to} & \begin{bmatrix} y_{i-1} & y_i \\ y_i & x^{b_i} \end{bmatrix} \succeq 0, \quad i=1 \dots q \\ & \begin{bmatrix} z & 1 \\ 1 & y_q \end{bmatrix} \succeq 0 \\ & x \in \mathcal{C} \end{array} \quad (553)$$

rather

$$x > 0, \quad z \geq x^{-\alpha} \Leftrightarrow \begin{bmatrix} y_{i-1} & y_i \\ y_i & x^{b_i} \\ z & 1 \\ 1 & y_q \end{bmatrix} \succeq 0, \quad i=1 \dots q \quad (554)$$

**3.3.1.3 positive inverted**

Now define vector  $t=[t_i, i=0 \dots q] \in \mathbb{R}^{q+1}$  with  $t_0=1$ . To implement an objective  $x^{1/\alpha}$  for quantized  $0 \leq \alpha < 1$  as in (548)

$$\begin{array}{ll} \text{minimize}_{x \in \mathbb{R}} & x^{1/\alpha} \\ \text{subject to} & x > 0 \\ & x \in \mathcal{C} \end{array} \equiv \begin{array}{ll} \text{minimize}_{x, y \in \mathbb{R}, t \in \mathbb{R}^{q+1}} & y \\ \text{subject to} & \begin{bmatrix} t_{i-1} & t_i \\ t_i & y^{b_i} \end{bmatrix} \succeq 0, \quad i=1 \dots q \\ & x = t_q > 0 \\ & x \in \mathcal{C} \end{array} \quad (555)$$

rather

$$x > 0, \quad y \geq x^{1/\alpha} \Leftrightarrow \begin{bmatrix} t_{i-1} & t_i \\ t_i & y^{b_i} \\ x = t_q & > 0 \end{bmatrix} \succeq 0, \quad i=1 \dots q \quad (556)$$

## 3.4 Affine function

A function  $f(X)$  is *affine* when it is continuous and has the dimensionally extensible form (confer §2.9.1.0.2)

$$f(X) = AX + B \quad (557)$$

All affine functions are simultaneously convex and concave. Both  $-AX + B$  and  $AX + B$ , for example, are convex functions of  $X$ . The linear functions constitute a proper subset of affine functions; *e.g.*, when  $B = \mathbf{0}$ , function  $f(X)$  is linear.

Unlike other convex functions, affine function convexity holds with respect to any dimensionally compatible proper cone substituted into convexity definition (502). All affine functions satisfy a membership relation, for some normal cone, like (511). Affine multidimensional functions are more easily recognized by existence of no multiplicative multivariate terms and no polynomial terms of degree higher than 1; *id est*, entries of the function are characterized by only linear combinations of argument entries plus constants.

$A^T X A + B^T B$  is affine in  $X$ , for example. Trace is an affine function; actually, a real linear function expressible as inner product  $f(X) = \langle A, X \rangle$  with matrix  $A$  being the Identity. The real affine function in Figure 77 illustrates hyperplanes, in its domain, constituting contours of equal function-value (*level sets* (562)).

#### 3.4.0.0.1 Example. Engineering control.

[450, §2.2]<sup>3.14</sup>

For  $X \in \mathbb{S}^M$  and matrices  $A, B, Q, R$  of any compatible dimensions, for example, the expression  $XAX$  is not affine in  $X$  whereas

$$g(X) = \begin{bmatrix} R & B^T X \\ XB & Q + A^T X + XA \end{bmatrix} \quad (558)$$

is an affine multidimensional function. Such a function is typical in engineering control. [63] [172]  $\square$

(*confer* Figure 18) Any single- or many-valued inverse of an affine function is affine.

#### 3.4.0.0.2 Example. Linear objective.

Consider minimization of a real affine function  $f(z) = a^T z + b$  over the convex feasible set  $\mathcal{C}$  in its domain  $\mathbb{R}^2$  illustrated in Figure 77. Since scalar  $b$  is fixed, the problem posed is the same as the convex optimization

$$\begin{array}{ll} \text{minimize}_z & a^T z \\ \text{subject to} & z \in \mathcal{C} \end{array} \quad (559)$$

whose objective of minimization is a real linear function. Were convex set  $\mathcal{C}$  polyhedral (§2.12), then this problem would be called a *linear program*. Were convex set  $\mathcal{C}$  an intersection with a positive semidefinite cone, then this problem would be called a *semidefinite program*.

There are two distinct ways to visualize this problem: one in the objective function's domain  $\mathbb{R}^2$ , the other including the ambient space of the objective function's range as in  $\begin{bmatrix} \mathbb{R}^2 \\ \mathbb{R} \end{bmatrix}$ . Both visualizations are illustrated in Figure 77. Visualization in the function domain is easier because of lower dimension and because

- level sets (562) of any affine function are affine. (§2.1.9)

In this circumstance, the level sets are parallel hyperplanes with respect to  $\mathbb{R}^2$ . One solves optimization problem (559) graphically by finding that hyperplane intersecting feasible set  $\mathcal{C}$  furthest right (in the direction of negative gradient  $-a$  (§3.6)).  $\square$

---

<sup>3.14</sup>The interpretation from this citation of  $\{X \in \mathbb{S}^M \mid g(X) \succeq 0\}$  as “an intersection between a linear subspace and the cone of positive semidefinite matrices” is incorrect. (See §2.9.1.0.2 for a similar example.) The conditions they state under which strong duality holds for semidefinite programming are conservative. (*confer* §4.2.3.0.1)

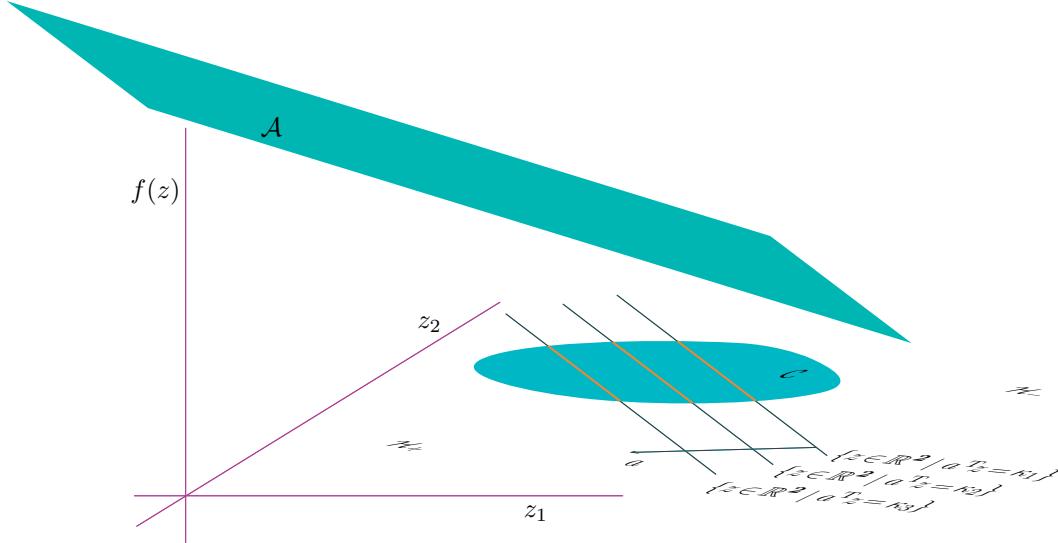


Figure 77: (confer Figure 31) Three hyperplanes intersecting convex set  $\mathcal{C} \subset \mathbb{R}^2$  from Figure 31. Cartesian axes in  $\mathbb{R}^3$ : Plotted is affine subset  $\mathcal{A} = f(\mathbb{R}^2) \subset \mathbb{R}^2 \times \mathbb{R}$ ; a plane with third dimension. We say sequence of hyperplanes, w.r.t domain  $\mathbb{R}^2$  of affine function  $f(z) = a^T z + b : \mathbb{R}^2 \rightarrow \mathbb{R}$ , is increasing in normal direction (Figure 29) because affine function increases in direction of gradient  $a$  (§3.6.0.0.3). Minimization of  $a^T z + b$  over  $\mathcal{C}$  is equivalent to minimization of  $a^T z$ .

When a differentiable convex objective function  $f$  is nonlinear, the negative gradient  $-\nabla f$  is a viable search direction (replacing  $-a$  in (559)). (§2.13.11.1, Figure 71) [173] Then the nonlinear objective function can be replaced with a dynamic linear objective; linear as in (559).

#### 3.4.0.0.3 Example. Support function.

[225, §C.2.1-§C.2.3.1]

For arbitrary set  $\mathcal{Y} \subseteq \mathbb{R}^n$ , its *support function*  $\sigma_{\mathcal{Y}}(a) : \mathbb{R}^n \rightarrow \mathbb{R}$  is defined

$$\sigma_{\mathcal{Y}}(a) \triangleq \sup_{z \in \mathcal{Y}} a^T z \quad (560)$$

a positively homogeneous function of direction  $a$  whose range contains  $\pm\infty$ . [280, p.135] For each  $z \in \mathcal{Y}$ ,  $a^T z$  is a linear function of vector  $a$ . Because  $\sigma_{\mathcal{Y}}(a)$  is a pointwise supremum of linear functions, it is convex in  $a$  (Figure 78). Application of the support function is illustrated in Figure 32a for one particular normal  $a$ . Given nonempty closed bounded convex sets  $\mathcal{Y}$  and  $\mathcal{Z}$  in  $\mathbb{R}^n$  and nonnegative scalars  $\beta$  and  $\gamma$  [434, p.234]

$$\sigma_{\beta\mathcal{Y}+\gamma\mathcal{Z}}(a) = \beta\sigma_{\mathcal{Y}}(a) + \gamma\sigma_{\mathcal{Z}}(a) \quad (561)$$

□

#### 3.4.0.0.4 Exercise. Level sets.

Given a function  $f$  and constant  $\kappa$ , its level sets are defined

$$\mathcal{L}_{\kappa} f \triangleq \{z \mid f(z) = \kappa\} \quad (562)$$

Give two distinct examples of convex function, that are not affine, having convex level sets. ▼

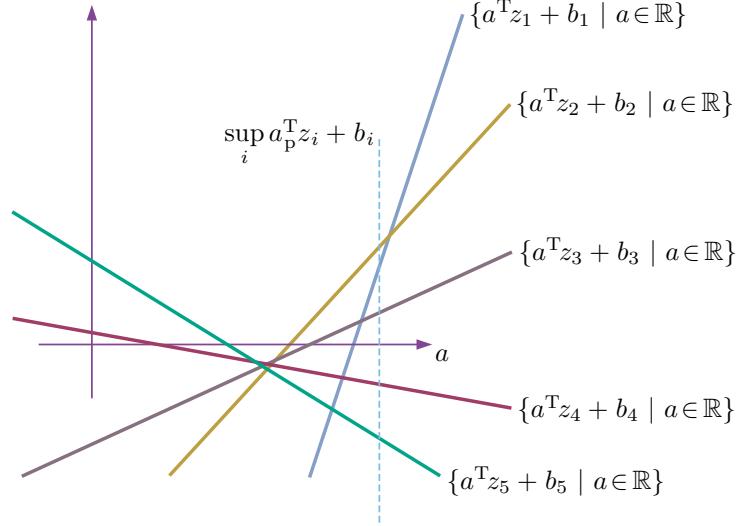


Figure 78: Pointwise supremum of any convex functions remains convex; by epigraph intersection. Supremum of affine functions in variable  $a$  evaluated at argument  $a_p$  is illustrated. Topmost affine function per  $a$  is supremum.

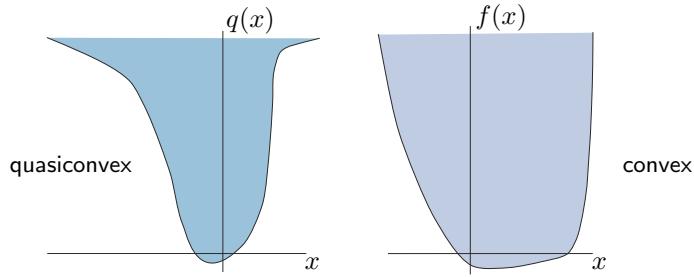


Figure 79: Quasiconvex function  $q$  epigraph is not necessarily convex, but convex function  $f$  epigraph is convex in any dimension. Sublevel sets are necessarily convex for either function, but sufficient only for quasiconvexity.

#### 3.4.0.0.5 Exercise. Epigraph intersection.

(confer Figure 78)

Draw three hyperplanes in  $\mathbb{R}^3$  representing  $\max(0, x) \triangleq \sup\{0, x_i | x \in \mathbb{R}^n\}$  in  $\mathbb{R}^2 \times \mathbb{R}$  to see why maximum of nonnegative vector entries is a convex function of variable  $x$ . What is the normal to each hyperplane?<sup>3.15</sup> Why is  $\max(x)$  convex? ▼

## 3.5 Epigraph, Sublevel set

It is well established that a continuous real function is convex if and only if its epigraph makes a convex set; [225] [343] [399] [434] [280] epigraph is the connection between convex sets and convex functions (p.169). Piecewise-continuous convex functions are admitted, thereby, and all invariant properties of convex sets carry over directly to convex functions. Generalization of *epigraph* to a vector-valued function  $f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}^M$  is straightforward: [329]

---

<sup>3.15</sup> Hint: p.198.

$$\text{epi } f \triangleq \{(X, t) \in \mathbb{R}^{p \times k} \times \mathbb{R}^M \mid X \in \text{dom } f, f(X) \underset{\mathbb{R}_+^M}{\preceq} t\} \quad (563)$$

*id est,*

$$f \text{ convex} \Leftrightarrow \text{epi } f \text{ convex} \quad (564)$$

Necessity is proven: [65, exer.3.60] Given any  $(X, u), (Y, v) \in \text{epi } f$ , we must show for all  $\mu \in [0, 1]$  that  $\mu(X, u) + (1-\mu)(Y, v) \in \text{epi } f$ ; *id est*, we must show

$$f(\mu X + (1-\mu)Y) \underset{\mathbb{R}_+^M}{\preceq} \mu u + (1-\mu)v \quad (565)$$

Yet this holds by definition because  $f(\mu X + (1-\mu)Y) \preceq \mu f(X) + (1-\mu)f(Y)$ .  
The converse also holds.  $\blacklozenge$

#### 3.5.0.0.1 Exercise. Epigraph sufficiency.

Prove that converse: Given any  $(X, u), (Y, v) \in \text{epi } f$ , if  $\mu(X, u) + (1-\mu)(Y, v) \in \text{epi } f$  holds for all  $\mu \in [0, 1]$ , then  $f$  must be convex.  $\blacktriangledown$

*Sublevel sets* of a convex real function are convex. Likewise, corresponding to each and every  $\nu \in \mathbb{R}^M$

$$\mathcal{L}_\nu f \triangleq \{X \in \text{dom } f \mid f(X) \underset{\mathbb{R}_+^M}{\preceq} \nu\} \subseteq \mathbb{R}^{p \times k} \quad (566)$$

sublevel sets of a convex vector-valued function are convex. As for convex real functions, the converse does not hold. (Figure 79)

To prove necessity of convex sublevel sets: For any  $X, Y \in \mathcal{L}_\nu f$  we must show for each and every  $\mu \in [0, 1]$  that  $\mu X + (1-\mu)Y \in \mathcal{L}_\nu f$ . By definition,

$$f(\mu X + (1-\mu)Y) \underset{\mathbb{R}_+^M}{\preceq} \mu f(X) + (1-\mu)f(Y) \underset{\mathbb{R}_+^M}{\preceq} \nu \quad (567)$$

$\blacklozenge$

When an epigraph (563) is artificially bounded above,  $t \preceq \nu$ , then the corresponding sublevel set can be regarded as an orthogonal projection of epigraph on the function domain.

Sense of the inequality is reversed in (563), for concave functions, and we use instead the nomenclature *hypograph*. Sense of the inequality in (566) is reversed, similarly, with each convex set then called *superlevel set*.

### 3.5.1 matrix pseudofractional function

Consider a real function of two variables

$$f(A, x) : \mathbb{S}_+^n \times \mathbb{R}^n \rightarrow \mathbb{R} = x^T A^\dagger x \quad (568)$$

on  $\text{dom } f = \mathbb{S}_+^n \times \mathcal{R}(A)$ . This function is convex simultaneously in both variables when variable matrix  $A$  belongs to the entire positive semidefinite cone  $\mathbb{S}_+^n$  and variable vector  $x$  is confined to range  $\mathcal{R}(A)$  of matrix  $A$ .

To explain this, we need only demonstrate that the function epigraph is convex.  
 Recall Schur-form (1663) from §A.4: for  $t \in \mathbb{R}$

$$\begin{aligned} G(A, z, t) &= \begin{bmatrix} A & z \\ z^T & t \end{bmatrix} \succeq 0 \\ &\Leftrightarrow \\ z^T(I - AA^\dagger) &= \mathbf{0} \\ t - z^T A^\dagger z &\geq 0 \\ A &\succeq 0 \end{aligned} \quad (569)$$

Inverse image of the positive semidefinite cone  $\mathbb{S}_+^{n+1}$  under affine mapping  $G(A, z, t)$  is convex by Theorem 2.1.9.0.1. Of the equivalent conditions for positive semidefiniteness of  $G$ , the first is an equality demanding that vector  $z$  belong to  $\mathcal{R}(A)$ . Function  $f(A, z) = z^T A^\dagger z$  is convex on convex domain  $\mathbb{S}_+^n \times \mathcal{R}(A)$  because the Cartesian product constituting its epigraph

$$\text{epi } f(A, z) = \{(A, z, t) \mid A \succeq 0, z \in \mathcal{R}(A), z^T A^\dagger z \leq t\} = G^{-1}(\mathbb{S}_+^{n+1}) \quad (570)$$

is convex.  $\blacklozenge$

### 3.5.1.0.1 Exercise. Matrix product function.

Continue §3.5.1 by introducing vector variable  $x$  and making the substitution  $z \leftarrow Ax$ . Because of matrix symmetry (§E), for all  $x \in \mathbb{R}^n$

$$f(A, z(x)) = z^T A^\dagger z = x^T A^T A^\dagger A x = x^T A x = f(A, x) \quad (571)$$

whose epigraph is

$$\text{epi } f(A, x) = \{(A, x, t) \mid A \succeq 0, x^T A x \leq t\} \quad (572)$$

Provide two simple explanations why  $f(A, x) = x^T A x$  is not a function convex simultaneously in positive semidefinite matrix  $A$  and vector  $x$  on  $\text{dom } f = \mathbb{S}_+^n \times \mathbb{R}^n$ .  $\blacktriangledown$

## 3.5.2 matrix fractional function

(confer §3.12.1) Now consider a real function of two variables on  $\text{dom } f = \mathbb{S}_+^n \times \mathbb{R}^n$  for small positive constant  $\epsilon$  (confer (2045))

$$f(A, x) = \epsilon x^T (A + \epsilon I)^{-1} x \quad (573)$$

where the inverse always exists by (1605). This function is convex simultaneously in both variables over the entire positive semidefinite cone  $\mathbb{S}_+^n$  and all  $x \in \mathbb{R}^n$ . This is explained:

Recall Schur-form (1666) from §A.4: for  $t \in \mathbb{R}$

$$\begin{aligned} G(A, z, t) &= \begin{bmatrix} A + \epsilon I & z \\ z^T & \epsilon^{-1} t \end{bmatrix} \succeq 0 \\ &\Leftrightarrow \\ t - \epsilon z^T (A + \epsilon I)^{-1} z &\geq 0 \\ A + \epsilon I &\succ 0 \end{aligned} \quad (574)$$

Inverse image of the positive semidefinite cone  $\mathbb{S}_+^{n+1}$  under affine mapping  $G(A, z, t)$  is convex by Theorem 2.1.9.0.1. Function  $f(A, z)$  is convex on  $\mathbb{S}_+^n \times \mathbb{R}^n$  because its epigraph is that inverse image:

$$\text{epi } f(A, z) = \{(A, z, t) \mid A + \epsilon I \succ 0, \epsilon z^T (A + \epsilon I)^{-1} z \leq t\} = G^{-1}(\mathbb{S}_+^{n+1}) \quad (575)$$

$\blacklozenge$

### 3.5.2.1 matrix fractional projector

Consider nonlinear function  $f$  having orthogonal projector  $W$  as argument:

$$f(W, x) = \epsilon x^T (W + \epsilon I)^{-1} x \quad (576)$$

Projection matrix  $W$  has property  $W^\dagger = W^T = W \succeq 0$  (2100). Any orthogonal projector can be decomposed into an outer product of orthonormal matrices  $W = UU^T$  where  $U^T U = I$  as explained in §E.3.2. From (2045) for any  $\epsilon > 0$  and idempotent symmetric  $W$ ,  $\epsilon(W + \epsilon I)^{-1} = I - (1 + \epsilon)^{-1}W$  from which

$$f(W, x) = \epsilon x^T (W + \epsilon I)^{-1} x = x^T (I - (1 + \epsilon)^{-1}W) x \quad (577)$$

$$\lim_{\epsilon \rightarrow 0^+} f(W, x) = \lim_{\epsilon \rightarrow 0^+} \epsilon x^T (W + \epsilon I)^{-1} x = x^T (I - W) x \quad (578)$$

where  $I - W$  is also an orthogonal projector (§E.2).

In §3.5.2 we learned that  $f(W, x) = \epsilon x^T (W + \epsilon I)^{-1} x$  is convex simultaneously in both variables over all  $x \in \mathbb{R}^n$  when  $W \in \mathbb{S}_+^n$  is confined to the entire positive semidefinite cone (including its boundary). It is now our goal to incorporate  $f$  into an optimization problem such that an optimal solution returned always comprises a projection matrix  $W$ . The set of orthogonal-projection matrices is a nonconvex subset of the positive semidefinite cone. So  $f$  cannot be convex on the projection matrices; its equivalent (for idempotent  $W$ )  $f(W, x) = x^T (I - (1 + \epsilon)^{-1}W) x$  cannot be convex simultaneously in both variables on either the positive semidefinite cone or orthogonal-projection matrices.

Suppose we allow  $\text{dom } f$  to constitute the entire positive semidefinite cone but confine  $W$  to a Fantope (91); e.g., for convex set  $\mathcal{C}$  and Fantope parametrized by  $0 < k < n$

$$\begin{aligned} & \underset{x \in \mathbb{R}^n, W \in \mathbb{S}^n}{\text{minimize}} && \epsilon x^T (W + \epsilon I)^{-1} x \\ & \text{subject to} && 0 \preceq W \preceq I \\ & && \text{tr } W = k \\ & && x \in \mathcal{C} \end{aligned} \quad (579)$$

Although this is a convex problem, there is no guarantee that optimal  $W$  is a projection matrix because only extreme points of a Fantope are orthogonal-projection matrices  $UU^T$ .

Let's try partitioning the problem into two convex parts (one for  $x$  and one for  $W$ ), substitute equivalence (577), and then iterate solution of convex quadratic problem

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && x^T (I - (1 + \epsilon)^{-1}W) x \\ & \text{subject to} && x \in \mathcal{C} \end{aligned} \quad (580)$$

with convex semidefinite problem

$$(a) \quad \begin{aligned} & \underset{W \in \mathbb{S}^n}{\text{minimize}} && x^{*\top} (I - (1 + \epsilon)^{-1}W) x^* \\ & \text{subject to} && 0 \preceq W \preceq I \\ & && \text{tr } W = k \end{aligned} \quad \equiv \quad \begin{aligned} & \underset{W \in \mathbb{S}^n}{\text{maximize}} && x^{*\top} W x^* \\ & \text{subject to} && 0 \preceq W \preceq I \\ & && \text{tr } W = k \end{aligned} \quad (581)$$

until convergence ( $x^*$  represents optimal solution to (580) at each successive iteration). The idea is to optimally solve for the partitioned variables which are later combined to solve the original problem (579). What makes this approach sound is that the constraints are separable, the partitioned feasible sets are not interdependent, and the fact that the original problem (though nonlinear) is convex simultaneously in both variables.<sup>3.16</sup>

---

<sup>3.16</sup>A convex problem has convex feasible set, and the objective *hypersurface* has one and only one global minimum. [338, p.123]

But partitioning alone does not guarantee a projector as solution. To make orthogonal projector  $W$  a certainty, we must invoke a known analytical solution to problem (581): Diagonalize optimal solution from problem (580)  $x^*x^{*\top} \triangleq Q\Lambda Q^\top$  (§A.5.1), and set  $U^* = Q(:, 1:k) \in \mathbb{R}^{n \times k}$  per (1872c);

$$W = U^*U^{*\top} = \frac{x^*x^{*\top}}{\|x^*\|^2} + Q(:, 2:k)Q(:, 2:k)^\top \quad (582)$$

Then optimal solution  $(W^*, x^*)$  to problem (579) is found, for small  $\epsilon$ , by iterating solution to problem (580) with projector solution (582) to convex problem (581).

**Proof.** Optimal vector  $x^*$  is orthogonal to the last  $n - 1$  columns of orthogonal matrix  $Q$ , so

$$f_{(580)}^* = \|x^*\|^2(1 - (1 + \epsilon)^{-1}) \quad (583)$$

after each iteration. Convergence of  $f_{(580)}^*$  is proven with the observation that iteration (580) (581a) is a nonincreasing sequence bounded below by 0. Any bounded monotonic sequence in  $\mathbb{R}$  is convergent. [289, §1.2] [43, §1.1] Expression (582) holds for projector  $W$  at each iteration, therefore  $\|x^*\|^2(1 - (1 + \epsilon)^{-1})$  must also represent the optimal objective value  $f_{(580)}^*$  at convergence.

Because the objective  $f_{(579)}$  from problem (579) is also bounded below by 0 on the same domain, this convergent optimal objective value  $f_{(580)}^*$  (for positive  $\epsilon$  arbitrarily close to 0) is sufficiently optimal for (579); *id est*, by (1849)

$$f_{(580)}^* \geq f_{(579)}^* \geq 0 \quad (584)$$

$$\lim_{\epsilon \rightarrow 0^+} f_{(580)}^* = 0 \quad (585)$$

Since optimal  $(U^*, x^*)$  from problem (580) is feasible to problem (579), and because their objectives are equivalent for projectors by (577), then converged  $(U^*, x^*)$  must also be optimal to (579) in the limit. Because problem (579) is convex, this represents a globally optimal solution. ♦

### 3.5.2.1.1 Exercise. Matrix fractional projector function class.

Show that there are larger positive values of  $\epsilon$  for which iteration (580) (581a) is equivalent to (579) and returns a projector  $W$ ; *id est*,  $\epsilon \rightarrow 0^+$  can be unnecessary. ▼

## 3.5.3 semidefinite program via Schur

*Schur complement* (1663) can be used to convert a projection problem to an optimization problem in *epigraph form*. Suppose, for example, we are presented with the constrained projection problem studied by Hayden & Wells in [210] (who provide analytical solution): Given  $A \in \mathbb{R}^{M \times M}$  and some full-rank matrix  $S \in \mathbb{R}^{M \times L}$  with  $L < M$

$$\begin{aligned} & \underset{X \in \mathbb{S}^M}{\text{minimize}} && \|A - X\|_F^2 \\ & \text{subject to} && S^T X S \succeq 0 \end{aligned} \quad (586)$$

Variable  $X$  is constrained to be positive semidefinite, but only on a subspace determined by  $S$ . First we write the epigraph form:

$$\begin{aligned} & \underset{X \in \mathbb{S}^M, t \in \mathbb{R}}{\text{minimize}} && t \\ & \text{subject to} && \|A - X\|_F^2 \leq t \\ & && S^T X S \succeq 0 \end{aligned} \quad (587)$$

Next we use Schur complement [309, §6.4.3] [278] and matrix vectorization (§2.2):

$$\begin{array}{ll} \underset{X \in \mathbb{S}^M, t \in \mathbb{R}}{\text{minimize}} & t \\ \text{subject to} & \begin{bmatrix} tI & \text{vec}(A - X) \\ \text{vec}(A - X)^T & 1 \end{bmatrix} \succeq 0 \\ & S^T X S \succeq 0 \end{array} \quad (588)$$

This semidefinite program (§4) is an epigraph form in disguise, equivalent to (586); it demonstrates how a quadratic objective or constraint can be converted to a semidefinite constraint.

Were problem (586) instead equivalently expressed without the norm square

$$\begin{array}{ll} \underset{X \in \mathbb{S}^M}{\text{minimize}} & \|A - X\|_F \\ \text{subject to} & S^T X S \succeq 0 \end{array} \quad (589)$$

then we get a subtle variation:

$$\begin{array}{ll} \underset{X \in \mathbb{S}^M, t \in \mathbb{R}}{\text{minimize}} & t \\ \text{subject to} & \|A - X\|_F \leq t \\ & S^T X S \succeq 0 \end{array} \quad (590)$$

that leads to an equivalent for (589) (and for (586) by (518))

$$\begin{array}{ll} \underset{X \in \mathbb{S}^M, t \in \mathbb{R}}{\text{minimize}} & t \\ \text{subject to} & \begin{bmatrix} tI & \text{vec}(A - X) \\ \text{vec}(A - X)^T & t \end{bmatrix} \succeq 0 \\ & S^T X S \succeq 0 \end{array} \quad (591)$$

### 3.5.3.0.1 Example. Schur anomaly.

Consider a problem, abstract in the convex constraint, given symmetric matrix  $A$

$$\begin{array}{ll} \underset{X \in \mathbb{S}^M}{\text{minimize}} & \|X\|_F^2 - \|A - X\|_F^2 \\ \text{subject to} & X \in \mathcal{C} \end{array} \quad (592)$$

the minimization of a difference of two quadratic functions each convex in matrix  $X$ . Observe equality

$$\|X\|_F^2 - \|A - X\|_F^2 = 2 \text{tr}(XA) - \|A\|_F^2 \quad (593)$$

So problem (592) is equivalent to the convex optimization

$$\begin{array}{ll} \underset{X \in \mathbb{S}^M}{\text{minimize}} & \text{tr}(XA) \\ \text{subject to} & X \in \mathcal{C} \end{array} \quad (594)$$

But this problem (592) does not have Schur-form

$$\begin{array}{ll} \underset{X \in \mathbb{S}^M, \alpha, t}{\text{minimize}} & t - \alpha \\ \text{subject to} & X \in \mathcal{C} \\ & \|X\|_F^2 \leq t \\ & \|A - X\|_F^2 \geq \alpha \end{array} \quad (595)$$

because the constraint in  $\alpha$  is nonconvex. (§2.9.1.0.1)  $\square$

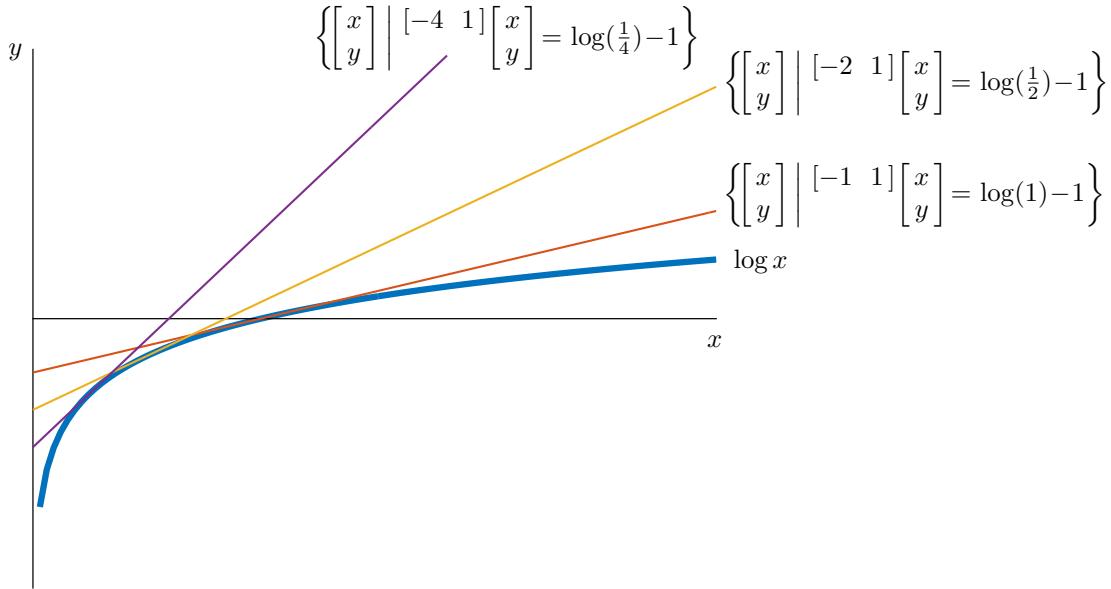


Figure 80: Three hyperplanes bounding hypograph of real function  $\log x$  from above.  $\log x \leq y$  for any  $[x \ y]^T$  belonging to a bounding hyperplane:  $\log x \leq x + \log(1) - 1$ ,  $\log x \leq 2x + \log(\frac{1}{2}) - 1$ ,  $\log x \leq 4x + \log(\frac{1}{4}) - 1$ .

Matrix 2-norm (*spectral norm*) coincides with largest singular value.

$$\|X\|_2 \triangleq \sup_{\|a\|=1} \|Xa\|_2 = \sigma(X)_1 = \sqrt{\lambda(X^T X)_1} = \begin{array}{ll} \text{minimize} & t \\ t \in \mathbb{R} & \\ \text{subject to} & \begin{bmatrix} tI & X \\ X^T & tI \end{bmatrix} \succeq 0 \end{array} \quad (596)$$

This supremum of a family of convex functions in  $X$  must be convex because it constitutes an intersection of epigraphs of convex functions.

### 3.5.4 Log

Suppose we want a variable and its logarithm to appear in constraints simultaneously. Such a problem formulation would generally be nonconvex. For example,

$$\begin{array}{ll} \text{minimize} & \alpha^T x \\ x \in \text{intr } \mathbb{R}_+^n & \\ \text{subject to} & Ax \leq b \\ & Cy \leq d \\ & y = \log x \end{array} \quad (597)$$

where  $\log x : \text{intr } \mathbb{R}_+^n \rightarrow \mathbb{R}^n$  operates on each entry individually. On the nonnegative real line, the  $\log$  function of real variable is concave having convex hypograph. Nonconvex problem (597) is solvable by approximating the hypograph with many bounding

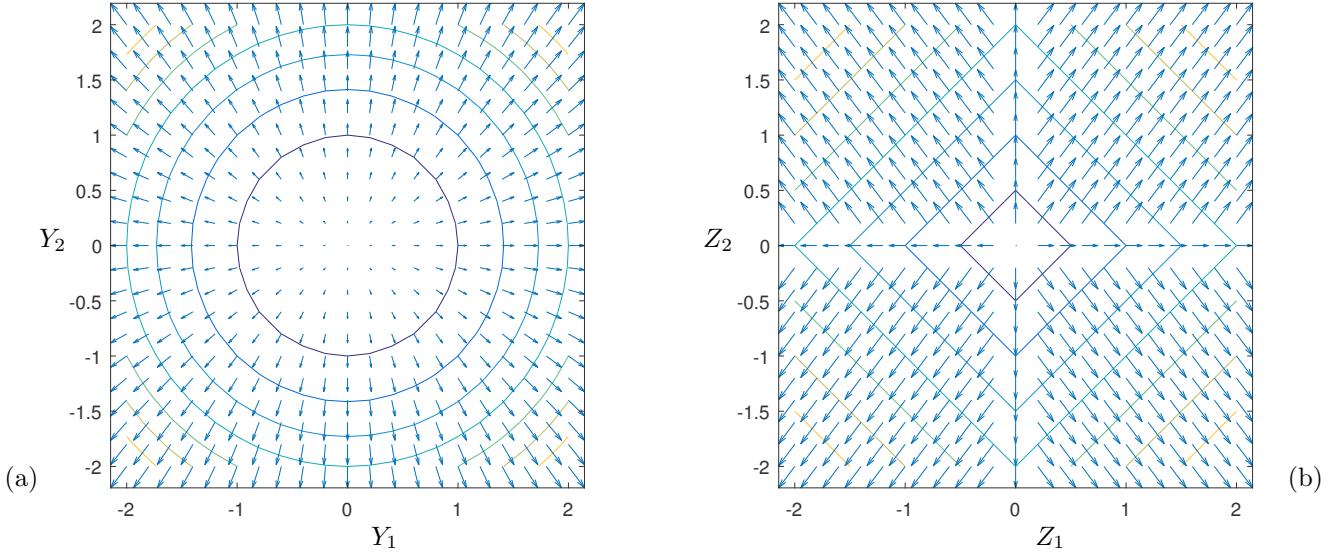


Figure 81: Contours are level sets; each defined by a constant function-value. Gradient in  $\mathbb{R}^2$  evaluated on grid over some open disc in domain of: (a) convex quadratic bowl  $f(Y) = Y^T Y : \mathbb{R}^2 \rightarrow \mathbb{R}$  illustrated in Figure 82 p.198, (b) 1-norm  $f(Z) = \|Z\|_1 : \mathbb{R}^2 \rightarrow \mathbb{R}$ .

hyperplanes as in Figure 80. In the problem, assignment  $y = \log x$  would be replaced with

$$\begin{bmatrix} \vdots & \vdots \\ -1 & 1 \\ -2 & 1 \\ -4 & 1 \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \preceq \begin{bmatrix} \vdots \\ \log(1)-1 \\ \log(\frac{1}{2})-1 \\ \log(\frac{1}{4})-1 \\ \vdots \end{bmatrix} \quad (598)$$

To succeed, this method requires nonnegative  $\alpha$  in the objective.

One can visualize this optimization in one dimension by imagining an objective function that pushes a feasible solution leftward along the  $x$  axis; driving toward the hypograph boundary which is the log function. In higher dimension, the same bounding hyperplane technique would be applied individually to each entry of  $\log$ . Accuracy to within any tolerance is ensured by increasing number of hyperplanes in vicinity of a solution.

## 3.6 Gradient

Gradient  $\nabla f$  of any differentiable multidimensional function  $f$  (formally defined in §D.1) maps each entry  $f_i$  to a space having the same dimension as the ambient space of its domain. Notation  $\nabla f$  is shorthand for gradient  $\nabla_x f(x)$  of  $f$  with respect to  $x$ .  $\nabla f(y)$  can mean  $\nabla_y f(y)$  or gradient  $\nabla_x f(y)$  of  $f(x)$  with respect to  $x$  evaluated at  $y$ ; a distinction that should become clear from context.

Gradient of a differentiable real function  $f(x) : \mathbb{R}^K \rightarrow \mathbb{R}$ , with respect to its vector argument, is uniquely defined

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_K} \end{bmatrix} \in \mathbb{R}^K \quad (1935)$$

while the second-order gradient of the twice differentiable real function, with respect to its vector argument, is traditionally called the *Hessian*<sup>3.17</sup>

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_K} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_K \partial x_1} & \frac{\partial^2 f(x)}{\partial x_K \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_K^2} \end{bmatrix} \in \mathbb{S}^K \quad (1936)$$

The gradient can be interpreted as a vector pointing in the direction of greatest change, [387, §15.6] or polar to direction of *steepest descent*<sup>3.18</sup>. The gradient can also be interpreted as that vector normal to a sublevel set; *e.g.*, Figure 83, Figure 71.

For the quadratic bowl in Figure 82, the gradient maps to  $\mathbb{R}^2$ ; illustrated in Figure 81. For a one-dimensional function of real variable, for example, the gradient evaluated at any point in the function domain is just the *slope* (or *derivative*) of that function there. (*confer* §D.1.4.1)

- For any differentiable multidimensional function, zero gradient  $\nabla f = \mathbf{0}$  is a condition necessary for its unconstrained minimization [173, §3.2]:

#### 3.6.0.0.1 Example. Projection on rank-1 subset.

For  $A \in \mathbb{S}^N$  having eigenvalues  $\lambda(A) = [\lambda_i] \in \mathbb{R}^N$ , consider the unconstrained nonconvex optimization that is a projection on the rank-1 subset (§2.9.2.1) of positive semidefinite cone  $\mathbb{S}_+^N$ : Defining  $\lambda_1 \triangleq \max_i \{\lambda(A)_i\}$  and corresponding eigenvector  $v_1$

$$\begin{aligned} \underset{x}{\text{minimize}} \quad & \|xx^T - A\|_F^2 = \underset{x}{\text{minimize}} \quad \text{tr}(xx^T(x^Tx) - 2Axx^T + A^TA) \\ &= \begin{cases} \|\lambda(A)\|^2, & \lambda_1 \leq 0 \\ \|\lambda(A)\|^2 - \lambda_1^2, & \lambda_1 > 0 \end{cases} \end{aligned} \quad (1866)$$

$$\arg \underset{x}{\text{minimize}} \quad \|xx^T - A\|_F^2 = \begin{cases} \mathbf{0}, & \lambda_1 \leq 0 \\ v_1 \sqrt{\lambda_1}, & \lambda_1 > 0 \end{cases} \quad (1867)$$

From (1965) and §D.2.1, the gradient of  $\|xx^T - A\|_F^2$  is

$$\nabla_x((x^Tx)^2 - 2x^TAx) = 4(x^Tx)x - 4Ax \quad (599)$$

Setting the gradient to  $\mathbf{0}$

$$Ax = x(x^Tx) \quad (600)$$

is necessary for optimal solution. Replace vector  $x$  with a normalized eigenvector  $v_i$  of  $A \in \mathbb{S}^N$ , corresponding to a positive eigenvalue  $\lambda_i$ , scaled by square root of that eigenvalue. Then (600) is satisfied

$$x \leftarrow v_i \sqrt{\lambda_i} \Rightarrow Av_i = v_i \lambda_i \quad (601)$$

<sup>3.17</sup>Jacobian is the Hessian transpose, so commonly confused in matrix calculus.

<sup>3.18</sup>Newton's direction  $-\nabla^2 f(x)^{-1} \nabla f(x)$  is better for optimization of nonlinear functions well approximated locally by a quadratic function. [173, p.105]

$xx^T = \lambda_i v_i v_i^T$  is a rank-1 matrix on the positive semidefinite cone boundary, and the minimum is achieved (§7.1.2) when  $\lambda_i = \lambda_1$  is the largest positive eigenvalue of  $A$ . If  $A$  has no positive eigenvalue, then  $x = \mathbf{0}$  yields the minimum.  $\square$

Differentiability is a prerequisite neither to convexity or to numerical solution of a convex optimization problem. The gradient provides a necessary and sufficient condition (355) (457) for optimality in the constrained case, nevertheless, as it does in the unconstrained case:

- For any differentiable multidimensional *convex* function, zero gradient  $\nabla f = \mathbf{0}$  is a necessary and sufficient condition for its unconstrained minimization [65, §5.5.3]:

#### 3.6.0.0.2 Example. Pseudoinverse $A^\dagger$ of matrix $A$ .

The pseudoinverse matrix is one particular solution from a convex set of solutions to an unconstrained convex optimization problem [181, §5.5.4]: given arbitrary  $A \in \mathbb{R}^{m \times n}$

$$\underset{X \in \mathbb{R}^{n \times m}}{\text{minimize}} \quad \|XA - I\|_F^2 \quad (602)$$

where

$$\|XA - I\|_F^2 = \text{tr}(A^T X^T X A - X A - A^T X^T + I) \quad (603)$$

whose gradient (§D.2.3)

$$\nabla_X \|XA - I\|_F^2 = 2(XAA^T - A^T) \quad (604)$$

vanishes when

$$XAA^T = A^T \quad (605)$$

We can make  $AA^T$  invertible by adding a positively scaled Identity: for any  $A \in \mathbb{R}^{m \times n}$

$$X \approx A^T(AA^T + tI)^{-1} \quad (606)$$

Invertibility is guaranteed for any finite positive value of  $t$  by (1605).

$$X = \lim_{t \rightarrow 0^+} (A^T A + tI)^{-1} A^T = \lim_{t \rightarrow 0^+} A^T (AA^T + tI)^{-1} \in \mathbb{R}^{n \times m} \quad (2046)$$

Then, in the limit  $t \rightarrow 0^+$ , matrix  $X$  becomes the pseudoinverse:  $X \rightarrow A^\dagger = X^*$ .

When matrix  $A$  is thin-or-square full-rank, in particular, then  $A^T A$  is invertible,  $(A^T A)^{-1} A^T = X^*$  is the pseudoinverse  $A^\dagger$ , and  $A^\dagger A = I$ . Starting with a minimization of  $\|AX - I\|_F^2$ , instead, invokes the second [*sic*] flavor in (2046): When matrix  $A$  is wide full-rank, then  $AA^T$  is invertible,  $A^T(AB^T)^{-1} = X^*$  becomes the pseudoinverse  $A^\dagger$ , and  $AA^\dagger = I$ . But (2046) always provides that unique pseudoinverse  $A^\dagger$ , regardless of shape or rank, that simultaneously minimizes  $\|AX - I\|_F^2$  and  $\|XA - I\|_F^2$ .  $\square$

#### 3.6.0.0.3 Example. Hyperplane, line, described by affine function.

Consider the real affine function of vector variable, (confer Figure 77)

$$f(x) : \mathbb{R}^p \rightarrow \mathbb{R} = a^T x + b \quad (607)$$

whose domain is  $\mathbb{R}^p$  and whose gradient  $\nabla f(x) = a$  is a vector constant (independent of  $x$ ). This function describes the real line  $\mathbb{R}$  (its range), and it describes a *nonvertical* [225, §B.1.2] hyperplane  $\partial\mathcal{H}$  in the space  $\mathbb{R}^p \times \mathbb{R}$  for any particular vector  $a$  (confer §2.4.2);

$$\partial\mathcal{H} = \left\{ \begin{bmatrix} x \\ a^T x + b \end{bmatrix} \mid x \in \mathbb{R}^p \right\} \subset \mathbb{R}^p \times \mathbb{R} \quad (608)$$

having nonzero normal

$$\eta = \begin{bmatrix} a \\ -1 \end{bmatrix} \in \mathbb{R}^p \times \mathbb{R} \quad (609)$$

This equivalence to a hyperplane holds only for real functions.<sup>3.19</sup> Epigraph of real affine function  $f(x)$  is therefore a halfspace in  $\begin{bmatrix} \mathbb{R}^p \\ \mathbb{R} \end{bmatrix}$ , so we have:

The real affine function is to convex functions  
as  
the hyperplane is to convex sets.

Similarly, the matrix-valued affine function of real variable  $x$ , for any particular matrix  $A \in \mathbb{R}^{M \times N}$

$$h(x) : \mathbb{R} \rightarrow \mathbb{R}^{M \times N} = Ax + B \quad (610)$$

describes a line in  $\mathbb{R}^{M \times N}$  in direction  $A$

$$\{Ax + B \mid x \in \mathbb{R}\} \subseteq \mathbb{R}^{M \times N} \quad (611)$$

and describes a line in  $\mathbb{R} \times \mathbb{R}^{M \times N}$

$$\left\{ \begin{bmatrix} x \\ Ax + B \end{bmatrix} \mid x \in \mathbb{R} \right\} \subset \mathbb{R} \times \mathbb{R}^{M \times N} \quad (612)$$

whose slope with respect to  $x$  is  $A$ .  $\square$

### 3.6.1 monotonic function

A real function of real argument is called *monotonic* when it is exclusively nonincreasing or nondecreasing over the whole of its domain. A real differentiable function of real argument is monotonic when its first derivative (not necessarily continuous) maintains sign over the function domain.

#### 3.6.1.0.1 Definition. Monotonicity.

In pointed closed convex cone  $\mathcal{K}$ , multidimensional function  $f(X)$  is

$$\begin{array}{ll} \text{nondecreasing monotonic when} & Y \succeq_{\mathcal{K}} X \Rightarrow f(Y) \succeq f(X) \\ \text{nonincreasing monotonic when} & Y \preceq_{\mathcal{K}} X \Rightarrow f(Y) \preceq f(X) \end{array} \quad (613)$$

$\forall X, Y \in \text{dom } f$ . Multidimensional function  $f(X)$  is

$$\begin{array}{ll} \text{increasing monotonic when} & Y \succ_{\mathcal{K}} X \Rightarrow f(Y) \succ f(X) \\ \text{decreasing monotonic when} & Y \prec_{\mathcal{K}} X \Rightarrow f(Y) \prec f(X) \end{array} \quad (614)$$

These latter inequalities define *strict monotonicity* when they hold over all  $X, Y \in \text{dom } f$ .  $\triangle$

---

<sup>3.19</sup>To prove that, consider a vector-valued affine function

$$f(x) : \mathbb{R}^p \rightarrow \mathbb{R}^M = Ax + b$$

having gradient  $\nabla f(x) = A^T \in \mathbb{R}^{p \times M}$ : The affine set

$$\left\{ \begin{bmatrix} x \\ Ax + b \end{bmatrix} \mid x \in \mathbb{R}^p \right\} \subset \mathbb{R}^p \times \mathbb{R}^M$$

is perpendicular to

$$\eta \triangleq \begin{bmatrix} \nabla f(x) \\ -I \end{bmatrix} \in \mathbb{R}^{p \times M} \times \mathbb{R}^{M \times M}$$

because

$$\eta^T \left( \begin{bmatrix} x \\ Ax + b \end{bmatrix} - \begin{bmatrix} 0 \\ b \end{bmatrix} \right) = 0 \quad \forall x \in \mathbb{R}^p$$

Yet  $\eta$  is a vector (in  $\mathbb{R}^p \times \mathbb{R}^M$ ) only when  $M = 1$ . ♦

For monotonicity of vector-valued functions,  $f$  compared with respect to the nonnegative orthant, it is necessary and sufficient for each entry  $f_i$  to be monotonic in the same sense.

Any affine function is monotonic. In  $\mathcal{K} = \mathbb{S}_+^M$ , for example,  $\text{tr}(Z^T X)$  is a nondecreasing monotonic function of matrix  $X \in \mathbb{S}^M$  when matrix constant  $Z$  is positive semidefinite; which follows from a result (384) of Fejér.

### 3.6.1.0.2 Exercise. Quasiconcave monotonic functions.

Prove:

$$A \succeq B \succeq 0 \Rightarrow \text{rank } A \geq \text{rank } B \quad (1645)$$

$$x \succeq y \succeq 0 \Rightarrow \text{card } x \geq \text{card } y \quad (615)$$

▼

A convex function can be characterized by another kind of nondecreasing monotonicity of its gradient:

### 3.6.1.0.3 Theorem. Gradient monotonicity. [225, §B.4.1.4] [58, §3.1 exer.20]

Given real differentiable function  $f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}$  with matrix argument on open convex domain, the condition

$$\langle \nabla f(Y) - \nabla f(X), Y - X \rangle \geq 0 \text{ for each and every } X, Y \in \text{dom } f \quad (616)$$

is necessary and sufficient for convexity of  $f$ . Strict inequality and *caveat*  $Y \neq X$  constitute necessary and sufficient conditions for strict convexity. ◇

### 3.6.1.0.4 Example. Composition of functions. [65, §3.2.4] [225, §B.2.1]

Monotonic functions play a vital role determining convexity of functions constructed by transformation. Given functions  $g : \mathbb{R}^k \rightarrow \mathbb{R}$  and  $h : \mathbb{R}^n \rightarrow \mathbb{R}^k$ , their composition  $f = g(h) : \mathbb{R}^n \rightarrow \mathbb{R}$  defined by

$$f(x) = g(h(x)), \quad \text{dom } f = \{x \in \text{dom } h \mid h(x) \in \text{dom } g\} \quad (617)$$

is convex if

- $g$  is convex nondecreasing monotonic **and**  $h$  is convex
- $g$  is convex nonincreasing monotonic **and**  $h$  is concave

is concave if

- $g$  is concave nondecreasing monotonic **and**  $h$  is concave
- $g$  is concave nonincreasing monotonic **and**  $h$  is convex

where  $\infty$  ( $-\infty$ ) is assigned to convex (concave)  $g$  when evaluated outside its domain. For differentiable functions, these rules are consequent to (1966).

- Convexity (concavity) of any  $g$  is preserved when  $h$  is affine. □

In particular, nondecreasing affine transformation of a convex (concave) function remains convex (concave). If  $f$  and  $g$  are nonnegative convex real functions, then  $(f(x)^k + g(x)^k)^{1/k}$  is also convex for integer  $k \geq 1$ . [279, p.44] A squared norm is convex having the same minimum because a squaring operation is convex nondecreasing monotonic on the nonnegative real line.

**3.6.1.0.5 Exercise.** *Anomalous composition.*

Composition of convex nondecreasing monotonic  $g = e^x$  with concave  $h = -x^2$ , each a real function, corresponds to no rule in Example 3.6.1.0.4. But  $g(h)$  is quasiconcave. (§3.14)<sup>3.20</sup> Does this kind of composition always produce quasiconcavity? ▼

**3.6.1.0.6 Exercise.** *Order of composition.*

Real function  $f = x^{-2}$  is convex on  $\mathbb{R}_+$  but not predicted so by results in Example 3.6.1.0.4 when  $g = h(x)^{-1}$  and  $h = x^2$ . Explain this anomaly. ▼

The following result, for product of real functions, is extensible to inner product of multidimensional functions on real domain:

**3.6.1.0.7 Exercise.** *Product and ratio of convex functions.* [65, exer.3.32]

In general the product or ratio of two convex functions is not convex. [258] However, there are some results that apply to functions on  $\mathbb{R}$  [real domain]. Prove the following.<sup>3.21</sup>

- (a) If  $f$  and  $g$  are convex, both nondecreasing (or nonincreasing), and positive functions on an interval, then  $fg$  is convex.
- (b) If  $f, g$  are concave, positive, with one nondecreasing and the other nonincreasing, then  $fg$  is concave.
- (c) If  $f$  is convex, nondecreasing, and positive, and  $g$  is concave, nonincreasing, and positive, then  $f/g$  is convex. ▼

## 3.7 First-order convexity condition, real function

Discretization of  $w \geq 0$  in (503) invites refocus to the real-valued function:

**3.7.0.0.1 Theorem.** *Necessary and sufficient convexity condition.* [43, §1.2] [342, §3] [65, §3.1.3] [150, §I.5.2] [366, §4.2] [453, §1.2.3] For  $f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}$  a real differentiable function with matrix argument on open convex domain, the condition (confer §D.1.7)

$$f(Y) \geq f(X) + \langle \nabla f(X), Y - X \rangle \quad \text{for each and every } X, Y \in \text{dom } f \quad (618)$$

is necessary and sufficient for convexity of  $f$ . Caveat  $Y \neq X$  and strict inequality again constitute necessary and sufficient conditions for strict convexity. [225, §B.4.1.1] ◇

When  $f(X) : \mathbb{R}^p \rightarrow \mathbb{R}$  is a real differentiable convex function with vector argument on open convex domain, there is simplification of the first-order condition (618); for each and every  $X, Y \in \text{dom } f$

$$f(Y) \geq f(X) + \nabla f(X)^T (Y - X) \quad (619)$$

From this we can find  $\underline{\partial H_-}$  a unique [434, p.220-229] nonvertical [225, §B.1.2] hyperplane (§2.4), expressed in terms of function gradient, supporting  $\text{epi } f$  at  $\begin{bmatrix} X \\ f(X) \end{bmatrix}$ : *videlicet*, defining  $f(Y \notin \text{dom } f) \triangleq \infty$  [65, §3.1.7]

$$\begin{bmatrix} Y \\ t \end{bmatrix} \in \text{epi } f \Leftrightarrow t \geq f(Y) \Rightarrow \begin{bmatrix} \nabla f(X)^T & -1 \end{bmatrix} \left( \begin{bmatrix} Y \\ t \end{bmatrix} - \begin{bmatrix} X \\ f(X) \end{bmatrix} \right) \leq 0 \quad (620)$$

<sup>3.20</sup>Hint: §3.15.

<sup>3.21</sup>Hint: Prove §3.6.1.0.7a by verifying Jensen's inequality ((502) at  $\mu = \frac{1}{2}$ ).

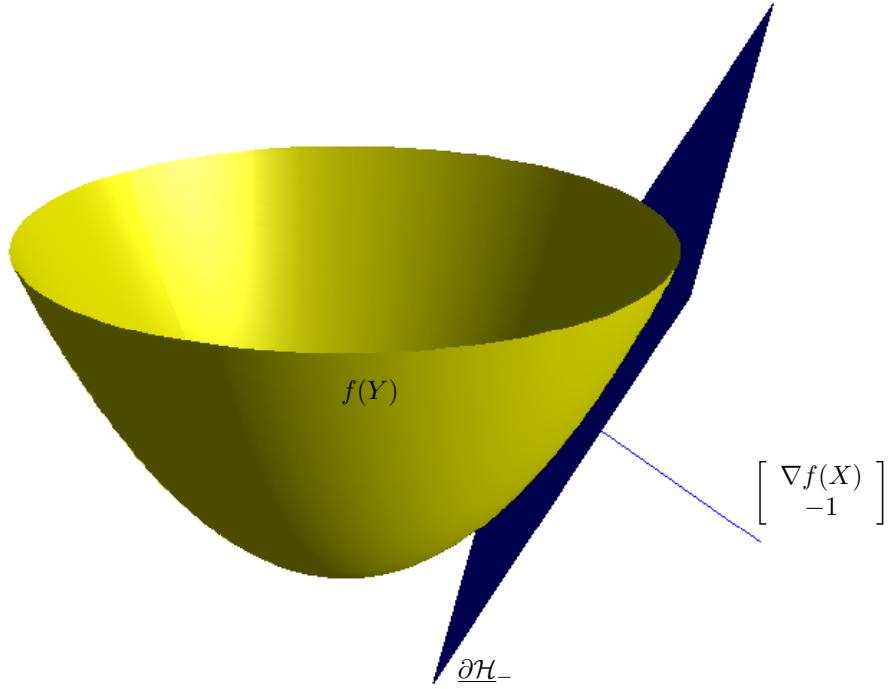


Figure 82: When a real function  $f$  is differentiable at each point in its open domain, there is an intuitive geometric interpretation of function convexity in terms of its gradient  $\nabla f$  (Figure 81 p.192) and its epigraph: Drawn is a convex quadratic bowl in  $\mathbb{R}^2 \times \mathbb{R}$  (*confer* Figure 186 p.551);  $f(Y) = Y^T Y : \mathbb{R}^2 \rightarrow \mathbb{R}$  versus  $Y$  on some open disc in  $\mathbb{R}^2$ . Unique strictly supporting hyperplane  $\underline{\partial}\mathcal{H}_- \subset \mathbb{R}^2 \times \mathbb{R}$  (only partially drawn) and its normal vector  $[\nabla f(X)^T \ -1]^T$ , at the particular point of support  $[X^T \ f(X)]^T$ , are illustrated. The interpretation: At each and every coordinate  $Y$ , normal  $[\nabla f(Y)^T \ -1]^T$  defines a unique hyperplane containing  $[Y^T \ f(Y)]^T$  and supporting the epigraph of convex differentiable  $f$ .

This means, for each and every point  $X$  in the domain of a convex real function  $f(X)$ , there exists a hyperplane  $\underline{\partial}\mathcal{H}_-$  in  $\mathbb{R}^p \times \mathbb{R}$  having normal  $\begin{bmatrix} \nabla f(X) \\ -1 \end{bmatrix}$  supporting the function epigraph at  $\begin{bmatrix} X \\ f(X) \end{bmatrix} \in \underline{\partial}\mathcal{H}_-$

$$\underline{\partial}\mathcal{H}_- = \left\{ \begin{bmatrix} Y \\ t \end{bmatrix} \in \begin{bmatrix} \mathbb{R}^p \\ \mathbb{R} \end{bmatrix} \mid [\nabla f(X)^T \ -1] \left( \begin{bmatrix} Y \\ t \end{bmatrix} - \begin{bmatrix} X \\ f(X) \end{bmatrix} \right) = 0 \right\} \quad (621)$$

Such a hyperplane is strictly supporting whenever a function is strictly convex. One such supporting hyperplane (*confer* Figure 32a) is illustrated in Figure 82 for a convex quadratic.

From (619) we deduce, for each and every  $X, Y \in \text{dom } f$  in the domain,

$$\nabla f(X)^T (Y - X) \geq 0 \Rightarrow f(Y) \geq f(X) \quad (622)$$

meaning, the gradient at  $X$  identifies a supporting hyperplane there in  $\mathbb{R}^p$

$$\{Y \in \mathbb{R}^p \mid \nabla f(X)^T (Y - X) = 0\} \quad (623)$$

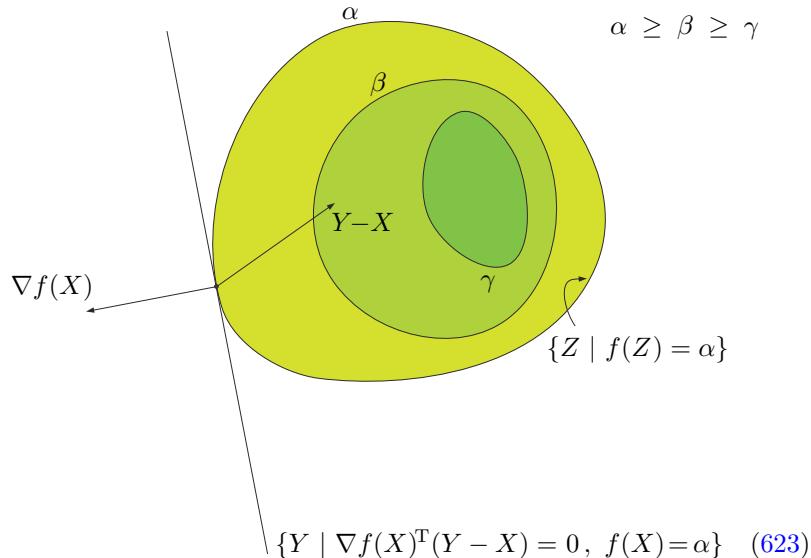


Figure 83: (confer Figure 71) Shown is a plausible contour plot in  $\mathbb{R}^2$  of some arbitrary real differentiable convex function  $f(Z)$  at selected levels  $\alpha$ ,  $\beta$ , and  $\gamma$ ; contours of equal level  $f$  (level sets) drawn in the function's domain. A convex function has convex sublevel sets  $\mathcal{L}_{f(X)}f$  (624). [343, §4.6] The sublevel set whose boundary is the level set at  $\alpha$ , for instance, comprises all shaded regions. For any particular convex function, the family comprising all its sublevel sets is nested. [225, p.75] Were sublevel sets not convex, we may certainly conclude the corresponding function is neither convex. Contour plots of real affine functions are illustrated in Figure 29 and Figure 77.

to the convex sublevel sets of convex function  $f$  (confer(566))

$$\mathcal{L}_{f(X)}f \triangleq \{Z \in \text{dom } f \mid f(Z) \leq f(X)\} \subseteq \mathbb{R}^p \quad (624)$$

illustrated for an arbitrary convex real function in Figure 83 and Figure 71. That supporting hyperplane is unique for twice differentiable  $f$ . [387, p.501]

### 3.7.1 Second-order convexity condition, real function

Disclosed in §3.9 and §3.13, a second-order condition for convexity of a real function corresponds to the one-dimensional case of a vector- or matrix-valued function.

## 3.8 First-order convexity condition, vector-valued $f$

Now consider the first-order necessary and sufficient condition for convexity of a vector-valued function: Differentiable function  $f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}^M$  is convex if and only if  $\text{dom } f$  is open, convex, and for each and every  $X, Y \in \text{dom } f$

$$f(Y) \succeq f(X) + \underset{\mathbb{R}_+^M}{\stackrel{\rightarrow Y-X}{df(X)}} \triangleq f(X) + \frac{d}{dt} \Big|_{t=0} f(X + t(Y - X)) \quad (625)$$

where  $\overset{\rightarrow Y-X}{df}(X)$  is the *directional derivative*<sup>3.22</sup> [387] [369] of  $f$  at  $X$  in direction  $Y-X$ . This, of course, follows from the real-valued function case: by dual generalized inequalities (§2.13.2.0.1),

$$f(Y) - f(X) - \overset{\rightarrow Y-X}{df}(X) \succeq 0 \Leftrightarrow \left\langle f(Y) - f(X) - \overset{\rightarrow Y-X}{df}(X), w \right\rangle \geq 0 \quad \forall w \succeq 0 \quad (626)$$

where

$$\overset{\rightarrow Y-X}{df}(X) = \begin{bmatrix} \text{tr}(\nabla f_1(X)^T(Y-X)) \\ \text{tr}(\nabla f_2(X)^T(Y-X)) \\ \vdots \\ \text{tr}(\nabla f_M(X)^T(Y-X)) \end{bmatrix} \in \mathbb{R}_+^M \quad (627)$$

Necessary and sufficient discretization (503) allows relaxation of the semiinfinite number of conditions  $\{w \succeq 0\}$  instead to  $\{w \in \{e_i, i=1 \dots M\}\}$  the extreme directions of the selfdual nonnegative orthant. Each extreme direction picks out a real entry  $f_i$  and  $\overset{\rightarrow Y-X}{df}(X)_i$  from the vector-valued function and its directional derivative, then Theorem 3.7.0.0.1 applies.

The vector-valued function case (625) is therefore a straightforward application of the first-order convexity condition for real functions to each entry of the vector-valued function.

### 3.9 Second-order convexity condition, vector-valued $f$

Again, by discretization (503), we are obliged only to consider each individual entry  $f_i$  of a vector-valued function  $f$ ; *id est*, the real functions  $\{f_i\}$ .

For  $f(X) : \mathbb{R}^p \rightarrow \mathbb{R}^M$ , a twice differentiable vector-valued function with vector argument on open convex domain,

$$\nabla^2 f_i(X) \succeq 0 \quad \forall X \in \text{dom } f, \quad i=1 \dots M \quad (628)$$

is necessary and sufficient for convexity of  $f$ . Condition (628) demands nonnegative curvature, intuitively, hence precluding points of inflection as in Figure 85 (p.207).

Obviously, when  $M=1$ , this convexity condition (628) also serves for a real function. Second-order convexity condition with matrix argument is deferred until §3.13.

Strict inequality in (628) provides only a sufficient condition ( $\Rightarrow$ ) for strict convexity, but that is nothing new; *videlicet*, strictly convex real function  $f_i(x)=x^4$  does not have positive second derivative at each and every  $x \in \mathbb{R}$ . Quadratic forms constitute a notable exception where the strict-case converse ( $\Leftarrow$ ) holds reliably.

#### 3.9.0.0.1 Example. Convex quadratic.

Real quadratic multivariate polynomial in matrix  $A$  and vector  $b$

$$x^T A x + 2b^T x + c \quad (629)$$

is convex if and only if  $A \succeq 0$ . Proof follows by observing second-order gradient: (§D.2.1)

$$\nabla_x^2(x^T A x + 2b^T x + c) = A + A^T \quad (630)$$

Because  $x^T(A + A^T)x = 2x^T A x$ , matrix  $A$  can be assumed symmetric.  $\square$

---

<sup>3.22</sup>We extend the traditional definition of directional derivative in §D.1.4 so that direction may be indicated by a vector or a matrix, thereby broadening scope of the Taylor series (§D.1.7). The right-hand side of inequality (625) is the first-order Taylor series expansion of  $f$  about  $X$ .

**3.9.0.0.2 Exercise.** *Real fractional function.* (confer §3.3, §3.5.2, §3.14.2.0.3)

Prove that real function  $f(x, y) = y/x$  is not convex on the first quadrant. Also exhibit this in a plot of the function. ( $f$  is *quasilinear* (p.208) on  $\{x > 0\}$  and nonmonotonic; even on the first quadrant.)  $\blacktriangledown$

**3.9.0.0.3 Exercise.** *One-dimensional stress function.*

Define  $|x - y| \triangleq \sqrt{(x - y)^2}$  and

$$X = [x_1 \cdots x_N] \in \mathbb{R}^{1 \times N} \quad (77)$$

Given symmetric nonnegative data  $[h_{ij}] \in \mathbb{S}^N \cap \mathbb{R}_+^{N \times N}$ , consider function

$$f(\text{vec } X) = \sum_{i=1}^{N-1} \sum_{j=i+1}^N (|x_i - x_j| - h_{ij})^2 \in \mathbb{R} \quad (1454)$$

Find a gradient and Hessian for  $f$ . Then explain why  $f$  is not a convex function; *id est*, why doesn't second-order condition (628) apply to the constant positive semidefinite Hessian matrix you found. For  $N=6$  and  $h_{ij}$  data from (1535), apply *line theorem* 3.13.0.0.1 to plot  $f$  along some arbitrary lines through its domain.  $\blacktriangledown$

**3.9.0.1 second-order  $\Rightarrow$  first-order condition**

For a twice-differentiable real function  $f_i(X) : \mathbb{R}^p \rightarrow \mathbb{R}$  having open domain, a consequence of the *mean value theorem* from calculus allows compression of its complete Taylor series expansion about  $X \in \text{dom } f_i$  (§D.1.7) to three terms: On some open interval of  $\|Y\|_2$ , so that each and every line segment  $[X, Y]$  belongs to  $\text{dom } f_i$ , there exists an  $\alpha \in [0, 1]$  such that [453, §1.2.3] [43, §1.1.4]

$$f_i(Y) = f_i(X) + \nabla f_i(X)^T(Y - X) + \frac{1}{2}(Y - X)^T \nabla^2 f_i(\alpha X + (1 - \alpha)Y)(Y - X) \quad (631)$$

The first-order condition for convexity (619) follows directly from this and the second-order condition (628).

## 3.10 Convex matrix-valued function

We need different tools for matrix argument: We are primarily interested in continuous matrix-valued functions  $g(X)$ . We choose symmetric  $g(X) \in \mathbb{S}^M$  because matrix-valued functions are most often compared (632) with respect to the positive semidefinite cone  $\mathbb{S}_+^M$  in the ambient space of symmetric matrices.  $\blacktriangledown$

**3.10.0.0.1 Definition.** *Convex matrix-valued function:***1) Matrix-definition.**

A function  $g(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{S}^M$  is convex in  $X$  iff  $\text{dom } g$  is a convex set and, for each and every  $Y, Z \in \text{dom } g$  and all  $0 \leq \mu \leq 1$  [242, §2.3.7]

$$g(\mu Y + (1 - \mu)Z) \underset{\mathbb{S}_+^M}{\preceq} \mu g(Y) + (1 - \mu)g(Z) \quad (632)$$

---

**3.23** Function symmetry is not a necessary requirement for convexity; indeed, for  $A \in \mathbb{R}^{m \times p}$  and  $B \in \mathbb{R}^{m \times k}$ ,  $g(X) = AX + B$  is a convex (affine) function in  $X$  on domain  $\mathbb{R}^{p \times k}$  with respect to the nonnegative orthant  $\mathbb{R}_+^{m \times k}$ . Symmetric convex functions share the same benefits as symmetric matrices. Horn & Johnson [228, §7.7] liken symmetric matrices to real numbers, and (symmetric) positive definite matrices to positive real numbers.

Reversing sense of the inequality flips this definition to concavity. Strict convexity is defined less a stroke of the pen in (632) similarly to (505).

**2) Scalar-definition.**

It follows that  $g(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{S}^M$  is convex in  $X$  iff  $w^T g(X)w : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}$  is convex in  $X$  for each and every  $\|w\| = 1$ ; shown by substituting the defining inequality (632). By dual generalized inequalities we have the equivalent but more broad criterion, (§2.13.6)

$$g \text{ convex w.r.t } \mathbb{S}_+^M \Leftrightarrow \langle W, g \rangle \text{ convex for each and every } W \succeq_{\mathbb{S}_+^M} 0 \quad (633)$$

Strict convexity on both sides requires *caveat*  $W \neq \mathbf{0}$ . Because the set of all extreme directions for the selfdual positive semidefinite cone (§2.9.2.7) comprises a minimal set of generators for that cone, discretization (§2.13.4.2.1) allows replacement of matrix  $W$  with symmetric dyad  $ww^T$  as proposed.  $\triangle$

**3.10.0.0.2 Example. Taxicab distance matrix.**

Consider an  $n$ -dimensional vector space  $\mathbf{R}^n$  with metric induced by the 1-norm. Then distance between points  $x_1$  and  $x_2$  is the norm of their difference:  $\|x_1 - x_2\|_1$ . Given a list of points arranged columnar in a matrix

$$X = [x_1 \cdots x_N] \in \mathbf{R}^{n \times N} \quad (77)$$

then we could define a taxicab distance matrix

$$\begin{aligned} \mathbf{D}_1(X) &\triangleq (I \otimes \mathbf{1}_n^T) |\text{vec}(X)\mathbf{1}^T - \mathbf{1} \otimes X| \in \mathbf{S}_h^N \cap \mathbf{R}_+^{N \times N} \\ &= \begin{bmatrix} 0 & \|x_1 - x_2\|_1 & \|x_1 - x_3\|_1 & \cdots & \|x_1 - x_N\|_1 \\ \|x_1 - x_2\|_1 & 0 & \|x_2 - x_3\|_1 & \cdots & \|x_2 - x_N\|_1 \\ \|x_1 - x_3\|_1 & \|x_2 - x_3\|_1 & 0 & \cdots & \|x_3 - x_N\|_1 \\ \vdots & \vdots & & \ddots & \vdots \\ \|x_1 - x_N\|_1 & \|x_2 - x_N\|_1 & \|x_3 - x_N\|_1 & \cdots & 0 \end{bmatrix} \end{aligned} \quad (634)$$

where  $\mathbf{1}_n$  is a vector of ones having  $\dim \mathbf{1}_n = n$  and where  $\otimes$  represents Kronecker product. This matrix-valued function is convex with respect to the nonnegative orthant since, for each and every  $Y, Z \in \mathbf{R}^{n \times N}$  and all  $0 \leq \mu \leq 1$

$$\mathbf{D}_1(\mu Y + (1 - \mu)Z) \underset{\mathbf{R}_+^{N \times N}}{\preceq} \mu \mathbf{D}_1(Y) + (1 - \mu) \mathbf{D}_1(Z) \quad (635)$$

$\square$

**3.10.0.0.3 Exercise. 1-norm distance matrix.**

The 1-norm is called *taxicab distance* because, to go from one point to another in a city by car, road distance is a sum of grid lengths as in Figure 84. Prove (635).  $\blacktriangledown$

**3.10.0.0.4 Exercise. Binary distance matrix.**

Euclidean distance-square between  $n$ -dimensional real vectors  $x_i \in \mathbb{R}^n$  is

$$\|x_i - x_j\|^2 = (x_i - x_j)^T(x_i - x_j) = x_i^T x_i + x_j^T x_j - 2x_i^T x_j \quad (1029)$$

One-bit signals are common in audio (Figure 1) and image processing. When encoding greyscale images,  $x_i$  can represent the  $i^{\text{th}}$  vectorized image with each entry corresponding to a particular *pixel* intensity ranging over dark to light. When black and white images are encoded,  $q_i$  represents a vectorized image with each entry either 0 or 1. Such binary

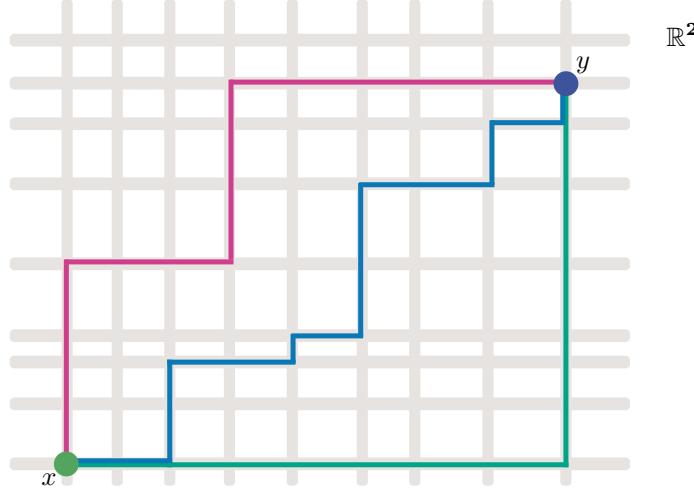


Figure 84: Path length, along any illustrated route from  $x$  to  $y$ , is identical to  $\|x - y\|_1$ . The same would hold in higher dimension, assuming no backtracking.

images are known as *halftone* (newspaper) or *fax* (*fac simile*). For  $n$ -dimensional binary vectors,  $s_i \in \mathbb{B}_\pm^n = \{-1, 1\}^n$  and  $q_i \in \mathbb{B}^n = \{0, 1\}^n$ , distance-square is expressed

$$\begin{aligned}\|s_i - s_j\|^2 &= (s_i - s_j)^T(s_i - s_j) = n + n - 2s_i^T s_j \\ \|q_i - q_j\|^2 &= (q_i - q_j)^T(q_i - q_j) = \mathbf{1}^T q_i + \mathbf{1}^T q_j - 2q_i^T q_j\end{aligned}\quad (636)$$

Each expression satisfies first properties of a metric (§5.2); importantly, nonnegativity and 0 selfdistance. Write a matrix expression, neater than (634), for a binary  $\mathbb{B}_\pm^n$  distance matrix. ▼

### 3.11 First-order convexity condition, matrix-valued $f$

From the *scalar-definition* (§3.10.0.0.1) of a convex matrix-valued function, for differentiable function  $g$  and for each and every real vector  $w$  of unit norm  $\|w\| = 1$ , we have

$$w^T g(Y) w \geq w^T g(X) w + w^T \overset{\rightarrow Y-X}{dg}(X) w \quad (637)$$

that follows immediately from the first-order condition (618) for convexity of a real function because

$$w^T \overset{\rightarrow Y-X}{dg}(X) w = \langle \nabla_X w^T g(X) w, Y - X \rangle \quad (638)$$

where  $\overset{\rightarrow Y-X}{dg}(X)$  is the directional derivative (§D.1.4) of function  $g$  at  $X$  in direction  $Y - X$ . By discretized dual generalized inequalities, (§2.13.6)

$$g(Y) - g(X) - \overset{\rightarrow Y-X}{dg}(X) \underset{\mathbb{S}_+^M}{\succeq} 0 \Leftrightarrow \left\langle g(Y) - g(X) - \overset{\rightarrow Y-X}{dg}(X), ww^T \right\rangle \underset{\mathbb{S}_+^{M^*}}{\geq} 0 \quad \forall ww^T(\succeq 0) \quad (639)$$

For each and every  $X, Y \in \text{dom } g$  (confer (625))

$$g(Y) \underset{\mathbb{S}_+^M}{\succeq} g(X) + \overset{\rightarrow Y-X}{dg}(X) \quad (640)$$

must therefore be necessary and sufficient for convexity of a matrix-valued function of matrix variable on open convex domain.

## 3.12 Epigraph of matrix-valued function, sublevel sets

We generalize *epigraph* to a continuous matrix-valued function  $g(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{S}^M$ : [35, p.155]

$$\text{epi } g \triangleq \{(X, T) \in \mathbb{R}^{p \times k} \times \mathbb{S}^M \mid X \in \text{dom } g, g(X) \preceq_{\mathbb{S}_+^M} T\} \quad (641)$$

from which it follows

$$g \text{ convex} \Leftrightarrow \text{epi } g \text{ convex} \quad (642)$$

Proof of necessity is similar to that in §3.5 on page 186.

Sublevel sets of a convex matrix-valued function corresponding to each and every  $S \in \mathbb{S}^M$  (*confer*(566))

$$\mathcal{L}_S g \triangleq \{X \in \text{dom } g \mid g(X) \preceq_{\mathbb{S}_+^M} S\} \subseteq \mathbb{R}^{p \times k} \quad (643)$$

are convex. There is no converse.

### 3.12.1 matrix fractional function *redux*

(*confer* §3.5.2) [35, p.155] Consider a matrix-valued function of two variables on  $\text{dom } g = \mathbb{S}_+^N \times \mathbb{R}^{n \times N}$  for small positive constant  $\epsilon$  (*confer*(2045))

$$g(A, X) = \epsilon X(A + \epsilon I)^{-1} X^T \quad (644)$$

where the inverse always exists by (1605). This function is convex simultaneously in both variables over the entire positive semidefinite cone  $\mathbb{S}_+^N$  and all  $X \in \mathbb{R}^{n \times N}$ . This is explained:

Recall Schur-form (1666) from §A.4: for  $T \in \mathbb{S}^n$

$$\begin{aligned} G(A, X, T) &= \begin{bmatrix} A + \epsilon I & X^T \\ X & \epsilon^{-1} T \end{bmatrix} \succeq 0 \\ &\Leftrightarrow \\ &T - \epsilon X(A + \epsilon I)^{-1} X^T \succeq 0 \\ &A + \epsilon I \succ 0 \end{aligned} \quad (645)$$

By Theorem 2.1.9.0.1, inverse image of the positive semidefinite cone  $\mathbb{S}_+^{N+n}$  under affine mapping  $G(A, X, T)$  is convex. Function  $g(A, X)$  is convex on  $\mathbb{S}_+^N \times \mathbb{R}^{n \times N}$  because its epigraph is that inverse image:

$$\begin{aligned} \text{epi } g(A, X) &= \{(A, X, T) \mid A + \epsilon I \succ 0, \epsilon X(A + \epsilon I)^{-1} X^T \preceq T\} \\ &= G^{-1}(\mathbb{S}_+^{N+n}) \end{aligned} \quad (646)$$

◆

## 3.13 Second-order convexity condition, matrix-valued $f$

The following *line theorem* is a potent tool for establishing convexity of a multidimensional function. To understand it, what is meant by *line* must first be solidified. Given a function  $g(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{S}^M$  and particular  $X, Y \in \mathbb{R}^{p \times k}$  not necessarily in that function's domain, then we say a line  $\{X + tY \mid t \in \mathbb{R}\}$  (infinite in extent) passes through  $\text{dom } g$  when  $X + tY \in \text{dom } g$  over some interval of  $t \in \mathbb{R}$ .

**3.13.0.0.1 Theorem.** *Line theorem.* (confer [65, §3.1.1])

Multidimensional function  $f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}^M$  or  $g(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{S}^M$  is convex in  $X$  if and only if it remains convex on the intersection of any line with its domain.  $\diamond$

Now we assume a twice differentiable function.

**3.13.0.0.2 Definition.** *Differentiable convex matrix-valued function.*

Matrix-valued function  $g(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{S}^M$  is convex in  $X$  iff  $\text{dom } g$  is an open convex set, and its second derivative  $g''(X + tY) : \mathbb{R} \rightarrow \mathbb{S}_+^M$  is positive semidefinite on each point of intersection along every line  $\{X + tY \mid t \in \mathbb{R}\}$  that intersects  $\text{dom } g$ ; *id est*, iff for each and every  $X, Y \in \mathbb{R}^{p \times k}$  such that  $X + tY \in \text{dom } g$  over some open interval of  $t \in \mathbb{R}$

$$\frac{d^2}{dt^2} g(X + tY) \succeq 0 \quad (647)$$

Similarly, if

$$\frac{d^2}{dt^2} g(X + tY) \succ 0 \quad (648)$$

then  $g$  is strictly convex; the converse is generally false. [65, §3.1.4]<sup>3.24</sup>  $\triangle$

**3.13.0.0.3 Example.** *Matrix inverse.*

(confer §3.3.1)

The matrix-valued function  $X^\mu$  is convex on  $\text{intr } \mathbb{S}_+^M$  for  $-1 \leq \mu \leq 0$  or  $1 \leq \mu \leq 2$  and concave for  $0 \leq \mu \leq 1$ . [65, §3.6.2] In particular, the function  $g(X) = X^{-1}$  is convex on  $\text{intr } \mathbb{S}_+^M$ . For each and every  $Y \in \mathbb{S}^M$  ([§D.2.1](#), [§A.3.1.0.5](#))

$$\frac{d^2}{dt^2} g(X + tY) = 2(X + tY)^{-1} Y (X + tY)^{-1} Y (X + tY)^{-1} \succeq 0 \quad (649)$$

on some open interval of  $t \in \mathbb{R}$  such that  $X + tY \succ 0$ . Hence,  $g(X)$  is convex in  $X$ . This result is extensible;<sup>3.25</sup>  $\text{tr } X^{-1}$  is convex on that same domain. [228, §7.6 prob.2] [[58](#), §3.1 exer.25]  $\square$

**3.13.0.0.4 Example.** *Matrix squared.*

Ionic real function  $f(x) = x^2$  is strictly convex on  $\mathbb{R}$ . The matrix-valued function  $g(X) = X^2$  is convex on the domain of symmetric matrices; for  $X, Y \in \mathbb{S}^M$  and any open interval of  $t \in \mathbb{R}$  ([§D.2.1](#))

$$\frac{d^2}{dt^2} g(X + tY) = \frac{d^2}{dt^2} (X + tY)^2 = 2Y^2 \quad (650)$$

which is positive semidefinite when  $Y$  is symmetric because then  $Y^2 = Y^T Y$  ([1611](#)).<sup>3.26</sup>

A more appropriate matrix-valued counterpart for  $f$  is  $g(X) = X^T X$  which is a convex function on domain  $\{X \in \mathbb{R}^{m \times n}\}$ , and strictly convex whenever  $X$  is thin-or-square full-rank. This matrix-valued function can be generalized to  $g(X) = X^T A X$  which is convex whenever matrix  $A$  is positive semidefinite ([p.560](#)), and strictly convex when  $A$  is positive definite and  $X$  is thin-or-square full-rank (Corollary [A.3.1.0.5](#)).  $\square$

<sup>3.24</sup>The strict-case converse is reliably true for quadratic forms.

<sup>3.25</sup> $d/dt \text{tr } g(X + tY) = \text{tr } d/dt g(X + tY)$ . [229, p.491]

<sup>3.26</sup>By ([1629](#)) in [§A.3.1](#), changing the domain instead to all symmetric and nonsymmetric positive semidefinite matrices, for example, will not produce a convex function.

**3.13.0.0.5 Exercise.** *Squared maps.*

Give seven examples of distinct polyhedra  $\mathcal{P}$  for which the set

$$\{X^T X \mid X \in \mathcal{P}\} \subseteq \mathbb{S}_+^n \quad (651)$$

were convex. Is this set convex, in general, for any polyhedron  $\mathcal{P}$ ? (confer (1379)(1386))  
Is the epigraph of function  $g(X) = X^T X$  convex for any polyhedral domain? ▼

**3.13.0.0.6 Exercise.** *Inverse square.*

(confer §3.12.1)

For positive scalar  $a$ , real function  $f(x) = ax^{-2}$  is convex on the nonnegative real line. Given positive definite matrix constant  $A$ , prove via *line theorem* that  $g(X) = \text{tr}((X^T A^{-1} X)^{-1})$  is generally not convex unless  $X \succ 0$ .<sup>3.27</sup> From this result, show how it follows via Definition 3.10.0.0.1-2 that  $h(X) = (X^T A^{-1} X)^{-1}$  is generally neither convex. ▼

**3.13.0.0.7 Example.** *Matrix exponential.*

The matrix-valued function  $g(X) = e^X : \mathbb{S}^M \rightarrow \mathbb{S}^M$  is convex on the subspace of *circulant* [192] symmetric matrices. Applying the *line theorem*, for all  $t \in \mathbb{R}$  and circulant  $X, Y \in \mathbb{S}^M$ , from Table D.2.7 we have

$$\frac{d^2}{dt^2} e^{X+tY} = Ye^{X+tY}Y \succeq 0, \quad (XY)^T = XY \quad (652)$$

because all circulant matrices are *commutative* and, for symmetric matrices,  $XY = YX \Leftrightarrow (XY)^T = XY$  (1628). Given symmetric argument, the matrix exponential always resides interior to the cone of positive semidefinite matrices in the symmetric matrix subspace;  $e^A \succ 0 \forall A \in \mathbb{S}^M$  (2043). Then for any matrix  $Y$  of compatible dimension,  $Y^T e^A Y$  is positive semidefinite. (§A.3.1.0.5)

The subspace of circulant symmetric matrices contains all diagonal matrices. The matrix exponential of any diagonal matrix  $e^A$  exponentiates each individual entry on the main diagonal. [281, §5.3] So, changing the function domain to the subspace of real diagonal matrices reduces the matrix exponential to a vector-valued function in an isometrically isomorphic subspace  $\mathbb{R}^M$ ; known convex (§3.1) from the real-valued function case [65, §3.1.5]. □

There are more methods for determining function convexity; [43] [65] [150] one can be more efficient than another depending on the function in hand.

**3.13.0.0.8 Exercise.**  *$\log \det$ .*

Matrix determinant is neither a convex or concave function, in general, but its inverse is convex when domain is restricted to interior of a positive semidefinite cone. [35, p.149] Show by three different methods: On interior of the positive semidefinite cone,  $\log \det X = -\log \det X^{-1}$  is concave. ▼

## 3.14 Quasiconvex

Quasiconvex functions [225] [366] [434] [273, §2] are valuable, pragmatically, because they are *unimodal* (by definition when nonmonotonic); a global minimum is guaranteed to exist over any convex set in the function domain; *e.g.*, Figure 85. That subset of the domain, corresponding to a global minimum, is convex. Optimal solution to quasiconvex problems is by method of bisection (**a.k.a. binary search**). [65, §4.2.5]

<sup>3.27</sup> Hint: §D.2.3.

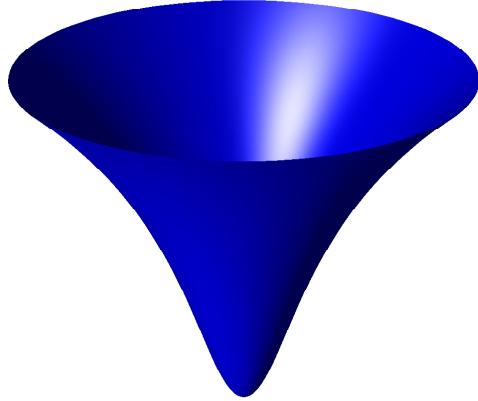


Figure 85: Iconic unimodal differentiable quasiconvex function of two variables graphed in  $\mathbb{R}^2 \times \mathbb{R}$  on some open disc in  $\mathbb{R}^2$ . Note reversal of curvature in direction of gradient. The negative of a quasiconvex function is quasiconcave, and *vice versa*.

**3.14.0.0.1 Definition.** *Quasiconvex function.* (confer(502))

$f(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{R}$  is a quasiconvex function of matrix  $X$  iff  $\text{dom } f$  is a convex set and for each and every  $Y, Z \in \text{dom } f$  and  $0 \leq \mu \leq 1$

$$f(\mu Y + (1 - \mu)Z) \leq \max\{f(Y), f(Z)\} \quad (653)$$

A quasiconcave function is similarly defined:

$$f(\mu Y + (1 - \mu)Z) \geq \min\{f(Y), f(Z)\} \quad (654)$$

*Caveat*  $Y \neq Z$  and strict inequality on an open interval  $0 < \mu < 1$  constitute necessary and sufficient conditions for strict quasiconvexity.  $\triangle$

Unlike convex functions, quasiconvex functions are not necessarily continuous; *e.g.*, quasiconcave  $\text{rank}(X)$  on  $\mathbb{S}_+^M$  ([§2.9.2.9.2](#)) and  $\text{card}(x)$  on  $\mathbb{R}_+^M$ . Although insufficient for convex functions, convexity of each and every sublevel set serves as a definition of quasiconvexity:

**3.14.0.0.2 Definition.** *Quasiconvex multidimensional function.*

Scalar-, vector-, or matrix-valued function  $g(X) : \mathbb{R}^{p \times k} \rightarrow \mathbb{S}^M$  is a quasiconvex function of matrix  $X$  iff  $\text{dom } g$  is a convex set and its sublevel set

$$\mathcal{L}_S g = \{X \in \text{dom } g \mid g(X) \preceq S\} \subseteq \mathbb{R}^{p \times k} \quad (643)$$

(corresponding to each and every  $S \in \mathbb{S}^M$ ) is convex. Vectors are compared with respect to the nonnegative orthant  $\mathbb{R}_+^M$  while matrices are with respect to the positive semidefinite cone  $\mathbb{S}_+^M$ .

Likewise, convexity of each and every *superlevel set*

$$\mathcal{L}^S g = \{X \in \text{dom } g \mid g(X) \succeq S\} \subseteq \mathbb{R}^{p \times k} \quad (655)$$

is necessary and sufficient for quasiconcavity.  $\triangle$

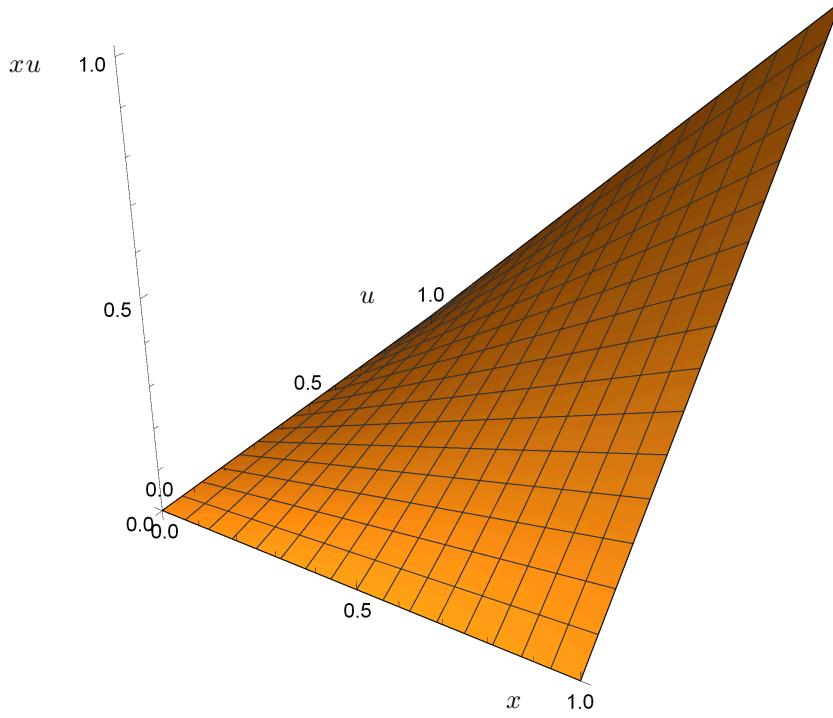


Figure 86: Quasiconcave strictly monotonic real product function  $xu$  is bowed (not affine) on the nonnegative orthants.

### 3.14.0.0.3 Exercise. Nonconvexity of matrix product.

Consider real function  $f$  on a positive definite domain

$$f(X) = \text{tr}(X_1 X_2), \quad X \triangleq \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \in \text{dom } f \triangleq \left[ \begin{array}{l} \text{rel intr } \mathbb{S}_+^N \\ \text{rel intr } \mathbb{S}_+^N \end{array} \right] \quad (656)$$

with superlevel sets

$$\mathcal{L}^s f = \{X \in \text{dom } f \mid \langle X_1, X_2 \rangle \geq s\} \quad (657)$$

Prove that  $f(X)$  is not quasiconcave except when  $N=1$ , nor is it quasiconvex unless  $X_1 = X_2$ .  $\blacktriangledown$

### 3.14.1 quasilinear

When a function is simultaneously quasiconvex and quasiconcave, it is called *quasilinear*. Quasilinear functions are completely determined by convex level sets. Multidimensional quasilinear functions are not necessarily monotonic; *e.g.*, Exercise 3.9.0.0.2.

One-dimensional functions  $x^3$  and  $e^x$  and vector-valued signum function  $\text{sgn}(x)$  are quasilinear, for example.

### 3.14.2 bilinear

Real bilinear (inner product) function<sup>3.28</sup>  $x^T u$  of vectors  $x$  and  $u$  is quasiconcave and strictly monotonic on the nonnegative orthants  $\mathbb{R}_+^\eta \times \mathbb{R}_+^\eta$  only when dimension  $\eta$  equals 1. (Figure 86)  $x^2 u$  and  $xu^2$  and biquadratic  $x^2 u^2$  are quasiconcave strictly monotonic, but over no more broad a domain.

When variable  $x \leftarrow \beta$  has dimension 1 but  $u$  is a vector of arbitrary dimension  $\eta$ , real function  $f(\beta, u) = \beta \mathbf{1}^T u = \beta \operatorname{tr} \delta(u)$  is quasiconcave strictly monotonic on the nonnegative orthants  $\mathbb{R}_+ \times \mathbb{R}_+^\eta$ .  $f(\beta, u)$  is quasiconcave strictly monotonic on  $\mathbb{R}_+ \times \mathbb{R}^\eta$  when  $\mathbf{1}^T u \geq 0$ .

#### 3.14.2.0.1 Proof. Domain of function

$$f(\beta, u) : \mathbb{R}_+ \times \mathbb{R}^\eta \rightarrow \mathbb{R}_+ = \beta \mathbf{1}^T u, \quad \mathbf{1}^T u \geq 0 \quad (658)$$

is a nonpointed polyhedral cone in  $\mathbb{R}^{\eta+1}$ , its range a halfline  $\mathbb{R}_+$ .

(quasiconcavity) Because this function spans an orthant, its 0-superlevel set is the deepest superlevel set and identical to its domain. Higher superlevel sets of the function, given some fixed nonzero scalar  $\zeta \geq 0$

$$\begin{aligned} \{\beta, u \mid f(\beta, u) \geq \zeta, \beta > 0, \mathbf{1}^T u \geq 0\} &= \{\beta, u \mid \beta \mathbf{1}^T u \geq \zeta, \beta > 0, \mathbf{1}^T u \geq 0\} \\ &= \{\beta, u \mid \mathbf{1}^T u \geq \frac{\zeta}{\beta}, \beta > 0\} \end{aligned} \quad (659)$$

are not polyhedral but they are convex because (§A.4)

$$\mathbf{1}^T u \geq \frac{\zeta}{\beta}, \quad \beta > 0 \quad \Leftrightarrow \quad \begin{bmatrix} \beta & \sqrt{\zeta} \\ \sqrt{\zeta} & \operatorname{tr} \delta(u) \end{bmatrix} \succeq 0 \quad (660)$$

and because inverse image of a positive semidefinite cone under affine transformation is convex by Theorem 2.1.9.0.1. Convex superlevel sets are necessary and sufficient for quasiconcavity by Definition 3.14.0.0.2.

(monotonicity) By Definition 3.6.1.0.1,

$$f \text{ is increasing monotonic when } \begin{bmatrix} \beta \\ u \end{bmatrix} \succ \begin{bmatrix} \tau \\ z \end{bmatrix} \Rightarrow f(\beta, u) \succ f(\tau, z) \quad (661)$$

for all  $\beta, u, \tau, z$  in the domain. Assuming  $\beta > \tau \geq 0$  and  $u \succ z$ , it follows that  $\mathbf{1}^T u > \mathbf{1}^T z$ . (Exercise 2.13.8.0.2) Therefore  $\beta \mathbf{1}^T u > \tau \mathbf{1}^T z$  and so  $f(\beta, u) = \beta \mathbf{1}^T u$  is strictly monotonic. ♦

#### 3.14.2.0.2 Example. Arbitrary magnitude analog filter design.

Analog filter design means determination of Laplace transfer function coefficients to meet specified tolerances in magnitude and phase over given frequencies. The problem posed by this example is to find a stable minimum phase filter  $H(j\omega)$ , of given order  $\eta$ , closest to specified samples of frequency-domain magnitude

$$\{g_i \mid \omega_i \in \Omega\} \quad (662)$$

We consider a recursive filter whose real causal impulse response has Laplace transform

$$H(s) = \frac{1 + b_1 s + b_2 s^2 + \dots + b_\eta s^\eta}{1 + a_1 s + a_2 s^2 + \dots + a_\eta s^\eta} \quad (663)$$

whose poles and zeros lie in the left halfplane ( $\operatorname{re} s = \sigma < 0$ ) and whose gain is unity at DC (at  $s = 0$ ). This transfer function is a ratio of polynomials in complex variable  $s$  defined

$$s \triangleq \sigma + j\omega \quad (664)$$

---

<sup>3.28</sup>Convex envelope of bilinear functions is well known. [4]

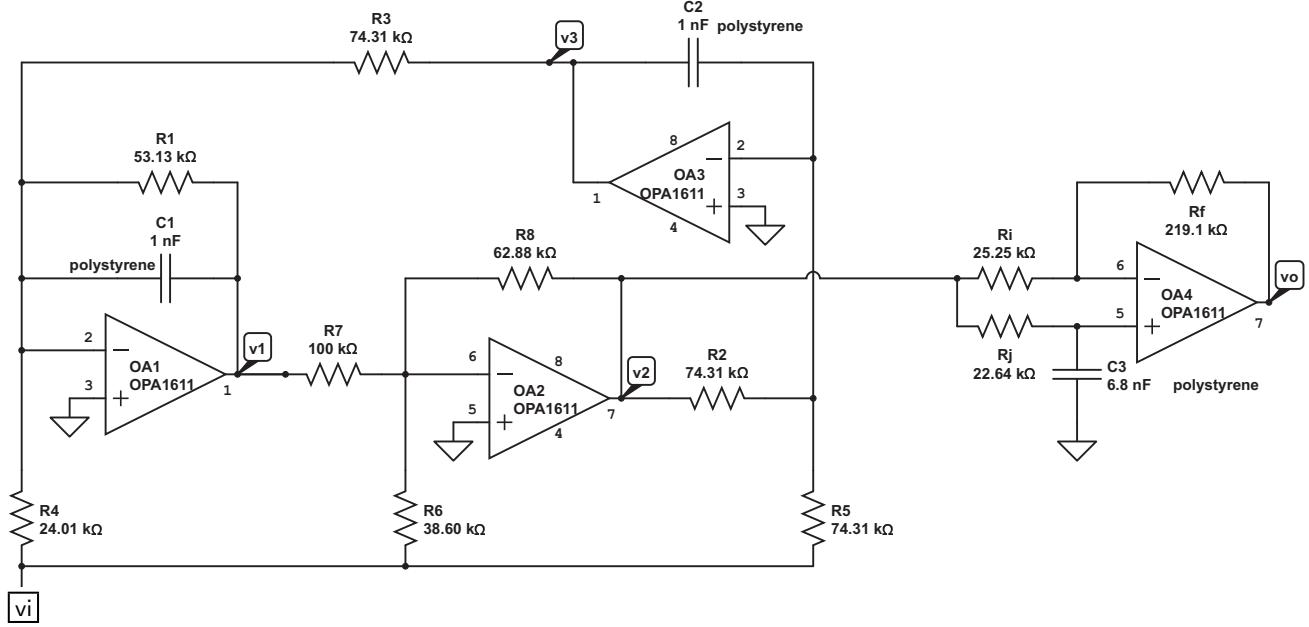


Figure 87: Third-order filter ( $\eta=3$ , confer Figure 90a). It is hard to introduce zeros of transmission, as we have succeeded in doing here [164] (Figure 88), because some *op amps* become unstable below unity gain by design.  $R_i$  and  $R_j$  minimize input offset voltage amplification.  $R_3 C_1 = R_2 C_2$  for minimum sensitivity.  $R_2 = R_5$  for unity gain at  $v_2$ .  $R_7(1/R_6 - R_2/(R_5 R_8)) = 1$  for unity gain at  $v_1$ . Bypass capacitors (6.8nF ceramic) reduce supply noise. polystyrene capacitors in signal path are essential for reducing distortion.

### realization

To reduce passive component sensitivity, physical implementation is facilitated by factoring Laplace transform (663) into parallel or cascade second-order sections which are needed to realize complex poles and zeros. Magnitude square of a second-order transfer function ( $\eta=2$ ) evaluated along the  $j\omega$  axis (the Fourier domain) is<sup>3.29</sup>

$$\frac{1 + v_1\omega^2 + v_2\omega^4}{1 + u_1\omega^2 + u_2\omega^4} = \frac{1 + (b_1^2 - 2b_2)\omega^2 + b_2^2\omega^4}{1 + (a_1^2 - 2a_2)\omega^2 + a_2^2\omega^4} \quad (665)$$

Coefficients  $b, a$  translate directly to passive component values. [164] Stability requires coefficients to obey  $0 < a_1$  and  $0 < a_2 \leq a_1^2/4$ , while minimum phase demands  $0 < b_1$  and  $0 < b_2 \leq b_1^2/4$ . These two requirements imply nonnegative magnitude square filter coefficients  $v, u$ .

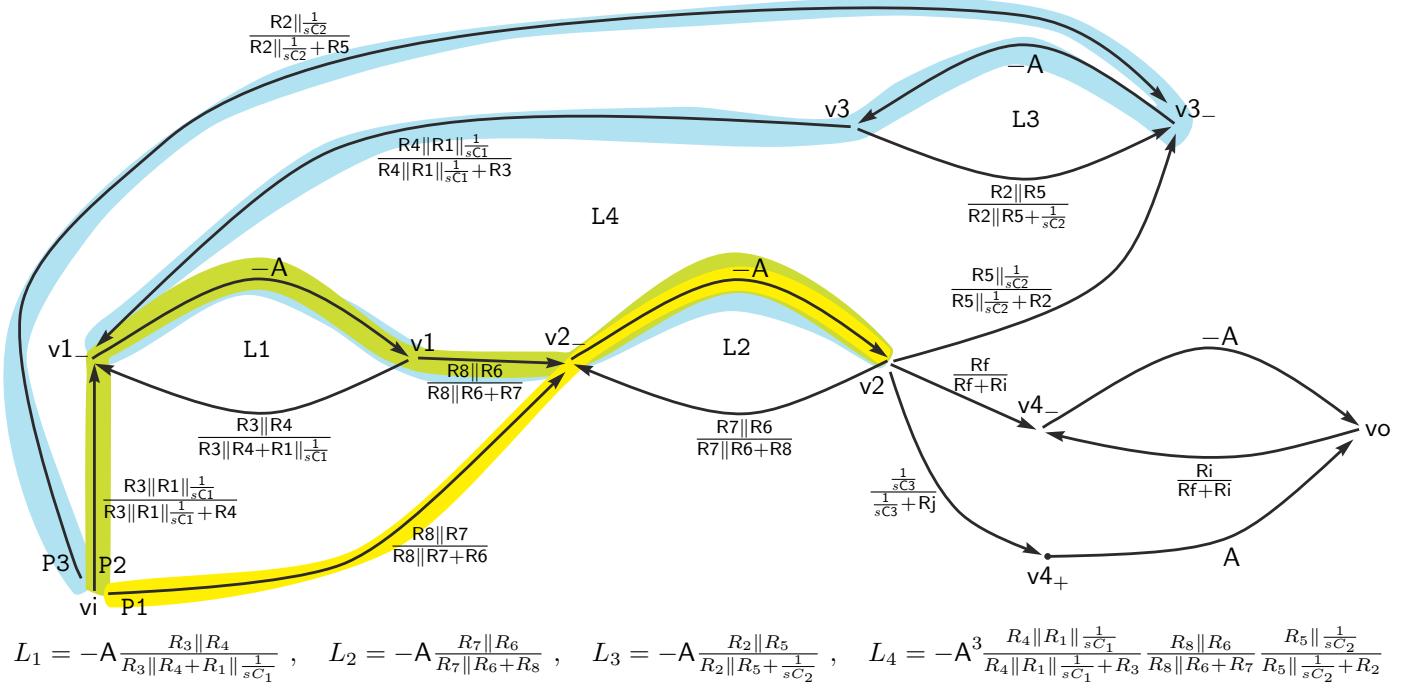
A cascade implementation of second-order sections can realize a high order unity gain filter. Magnitude square of  $\eta^{\text{th}}$  order transfer function  $H$ , evaluated along the  $j\omega$  axis, is

$$|H(j\omega)|^2 = H(j\omega)H(-j\omega) \triangleq \frac{V(\omega)}{U(\omega)} = \frac{1 + v_1\omega^2 + v_2\omega^4 + \dots + v_\eta\omega^{2\eta}}{1 + u_1\omega^2 + u_2\omega^4 + \dots + u_\eta\omega^{2\eta}} \quad (666)$$

A cascade of two second-order sections,  $\eta=4$  for example, has form

$$\frac{1 + \ddot{v}_1\omega^2 + \ddot{v}_2\omega^4}{1 + \ddot{u}_1\omega^2 + \ddot{u}_2\omega^4} \frac{1 + \ddot{v}_3\omega^2 + \ddot{v}_4\omega^4}{1 + \ddot{u}_3\omega^2 + \ddot{u}_4\omega^4} = \frac{1 + v_1\omega^2 + v_2\omega^4 + v_3\omega^6 + v_4\omega^8}{1 + u_1\omega^2 + u_2\omega^4 + u_3\omega^6 + u_4\omega^8} \quad (667)$$

<sup>3.29</sup>Real filter coefficient vectors  $b, a, v, u$  are independent of radian frequency  $\omega=2\pi f$ .



$$L_1 = -A \frac{R_3 \| R_4}{R_3 \| R_4 + R_1 \| \frac{1}{sC_1}} , \quad L_2 = -A \frac{R_7 \| R_6}{R_7 \| R_6 + R_8} , \quad L_3 = -A \frac{R_2 \| R_5}{R_2 \| R_5 + \frac{1}{sC_2}} , \quad L_4 = -A^3 \frac{R_4 \| R_1 \| \frac{1}{sC_1}}{R_4 \| R_1 \| \frac{1}{sC_1} + R_3} \frac{R_8 \| R_6}{R_8 \| R_6 + R_7} \frac{R_5 \| \frac{1}{sC_2}}{R_5 \| \frac{1}{sC_2} + R_2}$$

$$\Delta = 1 - L_1 - L_2 - L_3 - L_4 + L_1 L_2 + L_1 L_3 + L_2 L_3 - L_1 L_2 L_3$$

$$P_1 = -A \frac{R_8 \| R_7}{R_8 \| R_7 + R_6} , \quad P_2 = A^2 \frac{R_3 \| R_1 \| \frac{1}{sC_1}}{R_3 \| R_1 \| \frac{1}{sC_1} + R_4} \frac{R_8 \| R_6}{R_8 \| R_6 + R_7} , \quad P_3 = -A^3 \frac{R_2 \| \frac{1}{sC_2}}{R_2 \| \frac{1}{sC_2} + R_5} \frac{R_4 \| R_1 \| \frac{1}{sC_1}}{R_4 \| R_1 \| \frac{1}{sC_1} + R_3} \frac{R_8 \| R_6}{R_8 \| R_6 + R_7}$$

$$\Delta_1 = 1 - L_1 - L_3 + L_1 L_3 , \quad \Delta_2 = 1 - L_3 , \quad \Delta_3 = 1$$

$$\frac{v_2(s)}{v_i(s)} = \frac{P_1 \Delta_1 + P_2 \Delta_2 + P_3 \Delta_3}{\Delta} \approx \frac{-\frac{1}{R_5} - \left( \frac{R_7}{R_1 R_6} - \frac{1}{R_4} \right) R_3 C_2 s - \frac{1}{R_6} R_7 C_1 R_3 C_2 s^2}{\frac{1}{R_2} + \frac{R_7}{R_1 R_8} R_3 C_2 s + \frac{1}{R_8} R_7 C_1 R_3 C_2 s^2} , \quad \frac{v_o(s)}{v_2(s)} \approx \frac{1 - \frac{R_f}{R_i} R_j C_3 s}{1 + R_j C_3 s}$$

Figure 88: Matching the audio filter circuit comprising Figure 87, in Laplace variable  $s$  (664), Mason flowgraph is constructed solely by *voltage division*. Transfer function reduced algebraically by *Mathematica* code [44]. Op amp finite *open loop gain*  $A \approx 3E6$  makes transfer function inexact but a close approximation because high order terms  $A^4$  predominate. Magnitude transfer function  $\left| \frac{v_o(j\omega)}{v_i(j\omega)} \right|$  graphed in Figure 90a.

where (confer §4.5.1.2.5)

$$\begin{aligned} v_1 &= \ddot{v}_1 + \ddot{v}_3 , & u_1 &= \ddot{u}_1 + \ddot{u}_3 \\ v_2 &= \ddot{v}_2 + \ddot{v}_4 + \ddot{v}_1 \ddot{v}_3 , & u_2 &= \ddot{u}_2 + \ddot{u}_4 + \ddot{u}_1 \ddot{u}_3 \\ v_3 &= \ddot{v}_2 \ddot{v}_3 + \ddot{v}_1 \ddot{v}_4 , & u_3 &= \ddot{u}_2 \ddot{u}_3 + \ddot{u}_1 \ddot{u}_4 \\ v_4 &= \ddot{v}_2 \ddot{v}_4 , & u_4 &= \ddot{u}_2 \ddot{u}_4 \end{aligned} \quad (668)$$

Odd  $\eta$  is implemented by cascading one first-order section.

To guise filter design as an optimization problem, magnitude square filter coefficients

$$v \triangleq \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_\eta \end{bmatrix} \in \mathbb{R}^\eta , \quad u \triangleq \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_\eta \end{bmatrix} \in \mathbb{R}^\eta \quad (669)$$

become the variables. Given a set of frequencies  $\Omega$  at which finite  $g(\omega)$  is sampled, our notation reflects this role reversal:

$$\frac{V_i}{U_i} = \frac{V_i(v)}{U_i(u)} \triangleq \frac{1 + v_1\omega_i^2 + v_2\omega_i^4 + \dots + v_\eta\omega_i^{2\eta}}{1 + u_1\omega_i^2 + u_2\omega_i^4 + \dots + u_\eta\omega_i^{2\eta}}, \quad \omega_i \in \Omega \quad (670)$$

A better filter design may be obtained if stability and minimum phase requirements are ignored by the optimization.<sup>3.30</sup> Then coefficients  $v, u$  need not be nonnegatively constrained. A *least peak deviation filter* design problem may be expressed

$$\begin{aligned} & \underset{v, u, \beta}{\text{minimize}} \quad \beta \\ & \text{subject to} \quad \frac{1}{\beta} \leq \frac{U_i}{V_i} g_i^2 \leq \beta, \quad \omega_i \in \Omega \\ & \quad V_i \geq 0, \quad U_i \geq 0, \quad \omega_i \in \Omega \end{aligned} \quad (671)$$

which is a nonconvex optimization problem because of a ratio in variables  $v$  &  $u$ . Eliminating the ratio:

$$\begin{aligned} & \underset{v, u, \beta}{\text{minimize}} \quad \beta \\ & \text{subject to} \quad U_i g_i^2 \leq \beta V_i, \quad \omega_i \in \Omega \\ & \quad V_i g_i^{-2} \leq \beta U_i, \quad \omega_i \in \Omega \\ & \quad V_i \geq 0, \quad U_i \geq 0, \quad \omega_i \in \Omega \end{aligned} \quad (672)$$

This equivalent problem is likewise nonconvex because of new products in variables  $\beta$  &  $v$  and  $\beta$  &  $u$ . The feasible set in variables  $v, u, \beta$  is also nonconvex.

But problem (672) can be solved by applying the fact that each product  $\beta V_i(v), \beta U_i(u)$  increases monotonically. We may then rely on convexity of the feasible set described by linear constraints for fixed  $\beta$ . One strategy is to choose a sufficiently large positive  $\beta$  so that (672) is initially feasible; say  $\beta = \beta_o$ . Then halve  $\beta$  until (672) becomes infeasible; say  $\beta = \beta_\infty$ . At infeasibility, backtrack to a feasible  $\beta$ . Toward a perfect fit, in the sense that all  $g_i$  are collocated, optimal  $\beta$  approaches 1. This proposed iteration represents a binary search (**a.k.a bisection**) for minimum  $\beta$ , say  $\beta = \beta^*$ , that is assumed to lie in the interval  $[\beta_\infty, \beta_o]$ :

$$\begin{aligned} & \beta = \beta_o \\ & \beta_\infty = 0 \\ & \text{for } k=1, 2, \dots \text{until convergence} \{ \\ & \quad \text{if (672) feasible} \{ \\ & \quad \quad \beta_o = \beta_{k-1} \\ & \quad \text{else} \\ & \quad \quad \beta_\infty = \beta_{k-1} \\ & \quad \} \\ & \quad \beta = \beta_k = (\beta_o + \beta_\infty)/2 \\ & \} \end{aligned} \quad (673)$$

Convergence to a global optimum is certain, within any desired tolerance in absence of numerical error, characterized by an ever narrowing gap between  $\beta_\infty$  and  $\beta_o$ . At convergence of (673), the nonnegative  $\beta^*$  that minimizes (671) is found.

---

<sup>3.30</sup>Poles and zeros of  $H(j\omega)$  possess conjugate symmetry because of its real coefficients. Poles and zeros of  $H(j\omega)H(-j\omega)$  possess both conjugate and real symmetry. Stability and minimum phase may be imposed postoptimization by picking only those poles and zeros, respectively from  $H(j\omega)H(-j\omega)$ , that reside in the left half  $s$ -plane.

**Proof.** Problem (672) may be equivalently written

$$\begin{aligned} & \underset{v, u, \beta, t, r}{\text{minimize}} \quad \beta \\ & \text{subject to} \quad U_i g_i^2 \leq r, \quad \omega_i \in \Omega \\ & \quad r \leq \beta V_i, \quad \omega_i \in \Omega \\ & \quad V_i g_i^{-2} \leq t, \quad \omega_i \in \Omega \\ & \quad t \leq \beta U_i, \quad \omega_i \in \Omega \\ & \quad V_i \geq 0, \quad U_i \geq 0, \quad \omega_i \in \Omega \end{aligned} \tag{674}$$

where  $\beta, t, r \in \mathbb{R}_+$  are implicitly nonnegative. The feasible set is not convex because  $r$  is variable in  $r \leq \beta V_i$ ; this implicit union in  $\mathbb{R}^3$  (confer §2.1.9.0.1) of superlevel sets is nonconvex. The same holds for  $t \leq \beta U_i$ . By Schur complement (§A.4), quasiconcave strictly monotonic functions  $\beta V_i$  and  $\beta U_i$  (§3.14.2.0.1) may be decomposed

$$\begin{aligned} \beta V_i \geq r, \quad \beta > 0 & \Leftrightarrow \begin{bmatrix} V_i & \sqrt{r} \\ \sqrt{r} & \beta \end{bmatrix} \succeq 0 \\ \beta U_i \geq t, \quad \beta > 0 & \Leftrightarrow \begin{bmatrix} \beta & \sqrt{t} \\ \sqrt{t} & U_i \end{bmatrix} \succeq 0 \end{aligned} \tag{675}$$

By Theorem A.3.1.0.4, problem (674) is thereby equivalent to

$$\begin{aligned} & \underset{v, u, \beta, t, r}{\text{minimize}} \quad [0 \ 1 \ 0] \begin{bmatrix} V(v) & \sqrt{r} & 0 \\ \sqrt{r} & \beta & \sqrt{t} \\ 0 & \sqrt{t} & U(u) \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ & \text{subject to} \quad U_i g_i^2 \leq r, \quad \omega_i \in \Omega \\ & \quad V_i g_i^{-2} \leq t, \quad \omega_i \in \Omega \\ & \quad \begin{bmatrix} V_i & \sqrt{r} & 0 \\ \sqrt{r} & \beta & \sqrt{t} \\ 0 & \sqrt{t} & U_i \end{bmatrix} \succeq 0, \quad \omega_i \in \Omega \end{aligned} \tag{676}$$

This means minimizing  $\beta$  is like simultaneously minimizing functions  $\beta V(v)$  and  $\beta U(u)$  over a Cartesian subspace that is the intersection of their domains; namely, over a line containing the  $\beta$  axis in an increasing direction. Minimization on any line retains quasiconcavity [127], while minimization on any line in an increasing direction maintains strict monotonicity. But local minima cannot be precluded because minimizing a strictly monotonic function over a more general convex feasible set does not imply increasing direction (§3.6.1.0.1, Figure 86) and because, to begin with, the feasible set in variables  $v, u, \beta, t, r$  is nonconvex due to  $\sqrt{r}$  and  $\sqrt{t}$ .

One recourse, to preclude local minima, is to perform bisection (673) on  $\beta$  to find its globally optimal value  $\beta^*$  in problem (672):

$$\begin{aligned} & \text{find} \quad v, u \\ & \text{subject to} \quad U_i g_i^2 \leq \beta V_i, \quad \omega_i \in \Omega \\ & \quad V_i g_i^{-2} \leq \beta U_i, \quad \omega_i \in \Omega \\ & \quad V_i \geq 0, \quad U_i > 0, \quad \omega_i \in \Omega \end{aligned} \tag{677}$$

Now we explain how bisection succeeds. At each iteration, the feasible set in variables  $v, u$  is convex because all the inequalities are linear when  $\beta$  is fixed. A necessary and sufficient condition for bisection to find  $\beta^*$  requires feasibility for all  $\beta \geq \beta^*$  and infeasibility for all  $0 \leq \beta < \beta^*$ .  $\beta$  is bounded below by 0 but not bounded above (nor need it be). Feasibility is guaranteed by existence of a pointed convex cone  $\mathcal{K}$  described by the linear inequalities

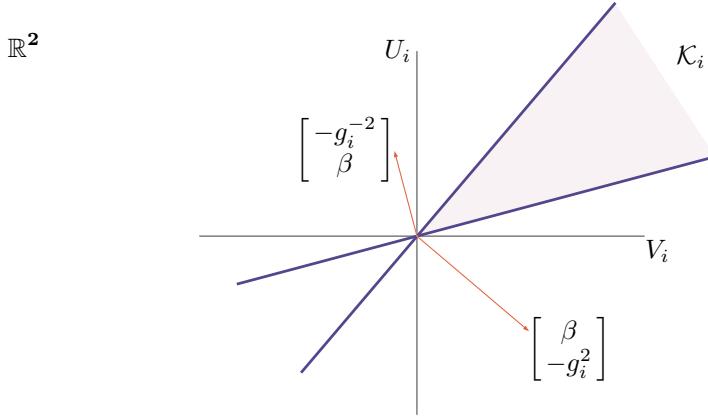


Figure 89: Linearity of bisection method for solution of a quasiconvex problem.

as an intersection of halfspaces about the origin in  $\mathbb{R}^2$ ; namely

$$\mathcal{K} = \left\{ \bigcap_i \mathcal{K}_i = \left\{ \begin{bmatrix} V_i \\ U_i \end{bmatrix} \in \mathbb{R}^2 \right\} \middle| \begin{array}{l} [\beta \ -g_i^2] \begin{bmatrix} V_i \\ U_i \end{bmatrix} \geq 0 \\ [-g_i^{-2} \ \beta] \begin{bmatrix} V_i \\ U_i \end{bmatrix} \geq 0 \\ \frac{V_i}{U_i} \geq 0 \\ \frac{V_i}{U_i} > 0 \end{array} \right\} \subseteq \mathbb{R}_+^2 \quad (678)$$

Cone  $\mathcal{K}$ , the feasible set in dependent variables  $V_i$  and  $U_i$ , ceases to exist for any  $\beta < \beta^*$ . This fact becomes evident by considering reach of the hyperplane normals which are linear in  $\beta$ :  $[-g_i^{-2} \ \beta]^T$  is confined to the second quadrant,  $[\beta \ -g_i^2]^T$  is confined to the fourth. As  $\beta \rightarrow \infty$ , the extreme directions of pointed convex cone  $\mathcal{K}_i$  approach the Cartesian axes in the first quadrant. As  $\beta \rightarrow 0$ , the extreme directions of  $\mathcal{K}_i$  collapse inward so as to narrow the cone then empty ( $\emptyset$ ) the feasible set. (Figure 89) When  $\beta = \beta^* \geq 1$ , cone intersection  $\mathcal{K}$  (§2.7.2.1) becomes an open ray emanating from the origin in the first quadrant.<sup>3.31</sup> Hyperplane-normal linearity in  $\beta$  and lack of an objective function here obviate the local-minima obstacle. ♦

A benefit, of this design methodology, is an unexpectedly low order  $\eta$  required to meet a given filter specification (662) to within reasonable tolerance  $\beta^*$ ; e.g., (Figure 87, Figure 88, Figure 90)

$$\begin{aligned} \frac{1}{2\pi}\Omega &= \{20 \ 30 \ 45 \ 60 \ 90 \ 125 \ 187 \ 250 \ 375 \ 500 \ 625 \ 750 \ 1000 \ 1250 \ 1500 \ 2000 \ 3000 \ 4000 \ 6000 \ 8000 \ 12000 \ 16000\} \\ 20 \log_{10}\{g_i\} &= \{0 \ 0 \ 0 \ 0 \ 0 \ 3 \ 6 \ 6 \ 9 \ 12 \ 12 \ 12 \ 6 \ 6 \ 6 \ 12 \ 21 \ 21 \ 21 \ 21 \ 21 \ 21 \ 21\} \end{aligned} \quad (679)$$

represents a particular individual's loss compensation, in dB, targeted by the earO assistive hearing device. Compensation levels are derived from an equal loudness hearing test then referenced to a *golden ear*.

<sup>3.31</sup> although no single  $\mathcal{K}_i$  is necessarily a ray.

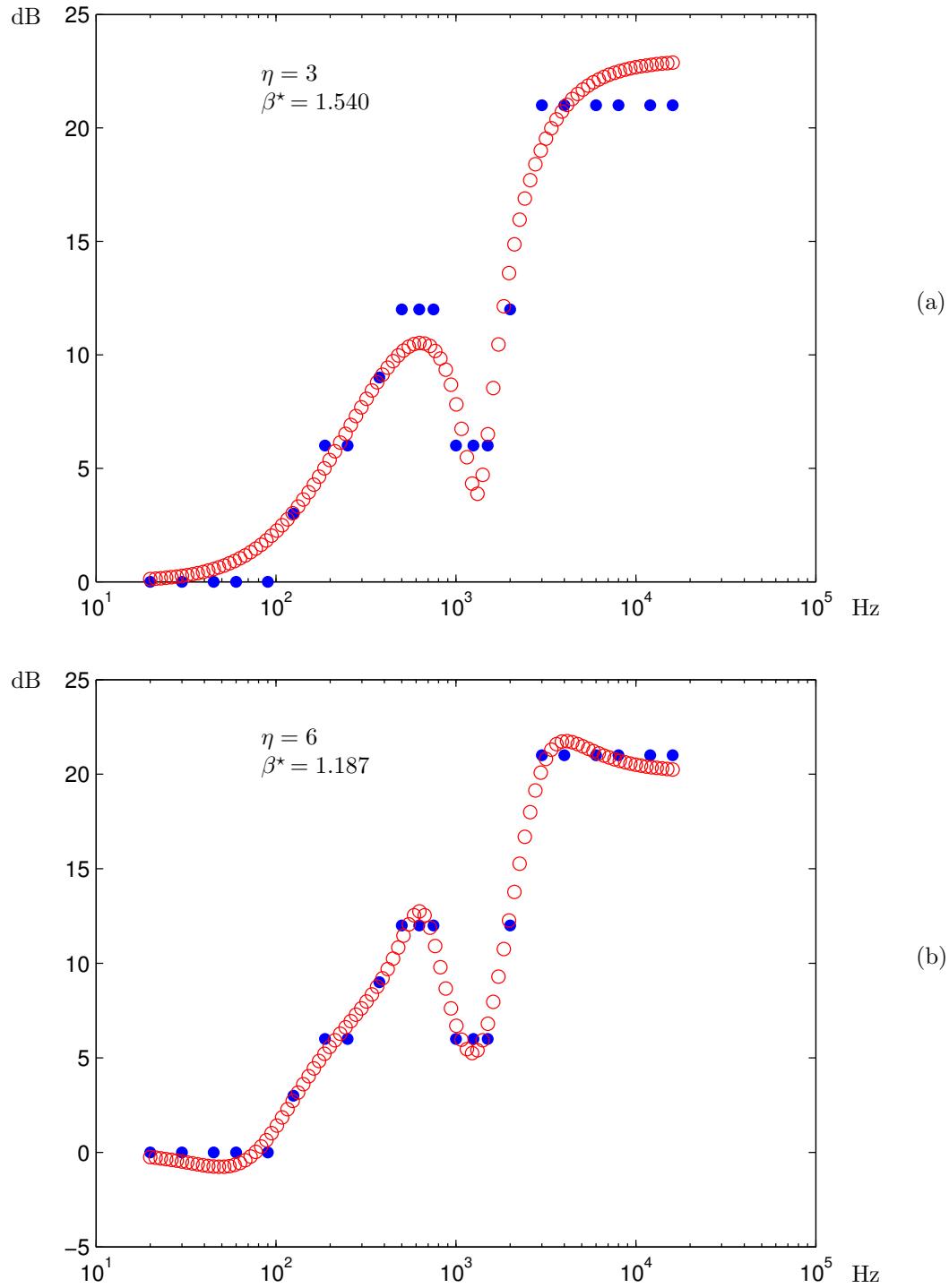


Figure 90: Arbitrary magnitude analog filter design specification represented by blue dots. Red circles represent fit: (a) third-order ( $\eta=3$ , confer Figure 87, Figure 88) and (b) sixth-order ( $\eta=6$ ). (Audiometrics for Glenna Mount.)

### precision

High powers of radian frequency, in magnitude square transfer function  $|H(j\omega)|^2$  (666), demand excessive precision from floating-point numerics. Even some second-order filter design optimization problems can cause the best linear program solvers to fail because double precision (64 bits, 52-bit mantissa) is inadequate. ([Saunders](#) provides a [quadruple precision solver](#): 128-bit wordlength, 112-bit mantissa.) High powers of frequency can be ameliorated simply by scaling specification (662):

$$\tilde{\Omega} \triangleq \frac{1}{c} \Omega \quad (680)$$

Scaled frequency-domain magnitude specification becomes

$$\{\tilde{g}_i \mid \tilde{\omega}_i \in \tilde{\Omega}\} = \{g_i \mid \omega_i \in \Omega\} \quad (681)$$

where

$$\tilde{\omega} \triangleq \frac{\omega}{c} \quad (682)$$

and  $c > 1$  is a constant. A filter is designed as before except magnitude square transfer function (666) is replaced with

$$|H(j\tilde{\omega})|^2 \triangleq \frac{\tilde{V}(\tilde{\omega})}{\tilde{U}(\tilde{\omega})} = \frac{1 + \tilde{v}_1\tilde{\omega}^2 + \tilde{v}_2\tilde{\omega}^4 + \dots + \tilde{v}_\eta\tilde{\omega}^{2\eta}}{1 + \tilde{u}_1\tilde{\omega}^2 + \tilde{u}_2\tilde{\omega}^4 + \dots + \tilde{u}_\eta\tilde{\omega}^{2\eta}} \quad (683)$$

Problem (672) is solved for magnitude square filter coefficients  $\tilde{v}, \tilde{u}$ :

$$\begin{aligned} \frac{\tilde{V}_i(\tilde{v})}{\tilde{U}_i(\tilde{u})} &\triangleq \frac{1 + \tilde{v}_1\tilde{\omega}_i^2 + \tilde{v}_2\tilde{\omega}_i^4 + \dots + \tilde{v}_\eta\tilde{\omega}_i^{2\eta}}{1 + \tilde{u}_1\tilde{\omega}_i^2 + \tilde{u}_2\tilde{\omega}_i^4 + \dots + \tilde{u}_\eta\tilde{\omega}_i^{2\eta}}, & \tilde{\omega}_i \in \tilde{\Omega} \end{aligned} \quad (684)$$

Optimal coefficients  $v, u$  are scaled replicas:

$$v^* \triangleq \begin{bmatrix} \frac{\tilde{v}_1^*}{c^2} \\ \frac{\tilde{v}_2^*}{c^4} \\ \vdots \\ \frac{\tilde{v}_\eta^*}{c^{2\eta}} \end{bmatrix}, \quad u^* \triangleq \begin{bmatrix} \frac{\tilde{u}_1^*}{c^2} \\ \frac{\tilde{u}_2^*}{c^4} \\ \vdots \\ \frac{\tilde{u}_\eta^*}{c^{2\eta}} \end{bmatrix} \quad (685)$$

The originally desired optimal fit is achieved as scaled frequency expands:

$$|H(j\omega)^*|^2 = \frac{1 + v_1^*(\tilde{\omega}c)^2 + v_2^*(\tilde{\omega}c)^4 + \dots + v_\eta^*(\tilde{\omega}c)^{2\eta}}{1 + u_1^*(\tilde{\omega}c)^2 + u_2^*(\tilde{\omega}c)^4 + \dots + u_\eta^*(\tilde{\omega}c)^{2\eta}} \quad (686)$$

□

#### 3.14.2.0.3 Exercise. Quasiconcave product function. (confer §3.9.0.0.2)

Show that vector-valued function  $\beta u : \mathbb{R}_+ \times \mathbb{R}_+^\eta \rightarrow \mathbb{R}_+^\eta$  is quasiconcave and strictly monotonic but not quasilinear. ▼

#### 3.14.2.0.4 Exercise. Sinusoid parameter estimation as quasiconvex problem.

Finding level and phase of a lone sinusoid is expressible as a convex problem. Show how successive estimation of frequency, alternating with approximation of level and phase, represents a method of solution to a *quasiconvex problem*; *id est*, minimization of a quasiconvex function over convex feasible set.<sup>3.32</sup> Is a condition on window length ( $M \rightarrow \infty$  or whole multiple of signal period) necessary for quasiconvexity? ▼

<sup>3.32</sup>This has been known for some time. [72, §3]

**3.14.2.0.5 Exercise.** *Digital filters of arbitrary magnitude frequency response.*

(confer §3.14.2.0.2) Like analog transfer function (663), a digital filter has transfer function

$$H(z) = \frac{b_0 + b_1 z^{-1} + b_2 z^{-2} + \dots + b_\eta z^{-\eta}}{a_0 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_\eta z^{-\eta}} \quad (687)$$

where  $z \triangleq e^{sT}$  where  $T$  is equal to sample period in seconds. Here  $H(z)$  is the  $z$  transform of a real causal discrete-time impulse response. [316] Magnitude square of a digital second-order transfer function,  $\eta=2$  for example, evaluated along the unit circle ( $z=e^{j\omega T}$ , the discrete-time Fourier domain) is

$$\frac{v_0 + v_1 \cos(\omega T) + v_2 \cos(2\omega T)}{u_0 + u_1 \cos(\omega T) + u_2 \cos(2\omega T)} = \frac{b_0^2 + b_1^2 + b_2^2 + 2b_1(b_0 + b_2) \cos(\omega T) + 2b_0 b_2 \cos(2\omega T)}{a_0^2 + a_1^2 + a_2^2 + 2a_1(a_0 + a_2) \cos(\omega T) + 2a_0 a_2 \cos(2\omega T)} \quad (688)$$

Evaluated at DC ( $z=1$ ) we get

$$\frac{v_0 + v_1 + v_2}{u_0 + u_1 + u_2} = \frac{(b_0 + b_1 + b_2)^2}{(a_0 + a_1 + a_2)^2} \quad (689)$$

More generally this means that if we require a unity gain filter at DC, then we acquire a new constraint

$$u_0 + u_1 + u_2 + \dots + u_\eta = v_0 + v_1 + v_2 + \dots + v_\eta = 1 \quad (690)$$

Magnitude square of  $\eta^{\text{th}}$  order digital transfer function  $H$ , evaluated along the unit circle

$$|H(e^{j\omega})|^2 = H(e^{j\omega})H(e^{-j\omega}) \triangleq \frac{V(\omega)}{U(\omega)} = \frac{v_0 + v_1 \cos(\omega T) + v_2 \cos(2\omega T) + \dots + v_\eta \cos(\eta\omega T)}{u_0 + u_1 \cos(\omega T) + u_2 \cos(2\omega T) + \dots + u_\eta \cos(\eta\omega T)} \quad (691)$$

no longer has high powers of radian frequency  $\omega$  as did its analog counterpart (666). A question naturally arises as to whether this digital magnitude square transfer function is more amenable to numerical computation. To answer this, choose a sample rate ( $F_s=1/T$ ) that is at least twice the highest analog frequency present in design specification (679);  $16000T \leq \frac{1}{2}$  must hold. Redefine

$$\frac{V_i}{U_i} = \frac{V_i(v)}{U_i(u)} \triangleq \frac{v_0 + v_1 \cos(\omega_i T) + v_2 \cos(2\omega_i T) + \dots + v_\eta \cos(\eta\omega_i T)}{u_0 + u_1 \cos(\omega_i T) + u_2 \cos(2\omega_i T) + \dots + u_\eta \cos(\eta\omega_i T)}, \quad \omega_i \in \Omega \quad (692)$$

which replaces (670). Solve problem (672), as before, but with new constraint (690):

$$\begin{aligned} & \underset{v, u, \beta}{\text{minimize}} \quad \beta \\ & \text{subject to} \quad U_i g_i^2 \leq \beta V_i, \quad \omega_i \in \Omega \\ & \quad V_i g_i^{-2} \leq \beta U_i, \quad \omega_i \in \Omega \\ & \quad V_i \geq 0, \quad U_i \geq 0, \quad \omega_i \in \Omega \\ & \quad U_0 = V_0 = 1 \end{aligned} \quad (693)$$

where the 0 subscript denotes phantom frequency  $\omega_0 \triangleq 0$  not present in specification (679). Assumption of unity gain at DC means corresponding implicit gain  $g_0$  is collocated and has value 1; *id est*,  $20 \log_{10} g_0 = 0$  dB. Optimal filter coefficients are very large for this particular design specification (679). Does numerical solution depend on sample rate; *e.g.*, does  $1/T=64000$  work better than  $1/T=32000$ ? Do solver numerics behave better here than for the corresponding analog filter design? Is higher filter order  $\eta$  consequently achievable? If in the affirmative, then one might accomplish a difficult analog design by first designing in the digital domain using a warped frequency specification obtained via bilinear transformation [316, §7.1.2].  $\blacktriangledown$

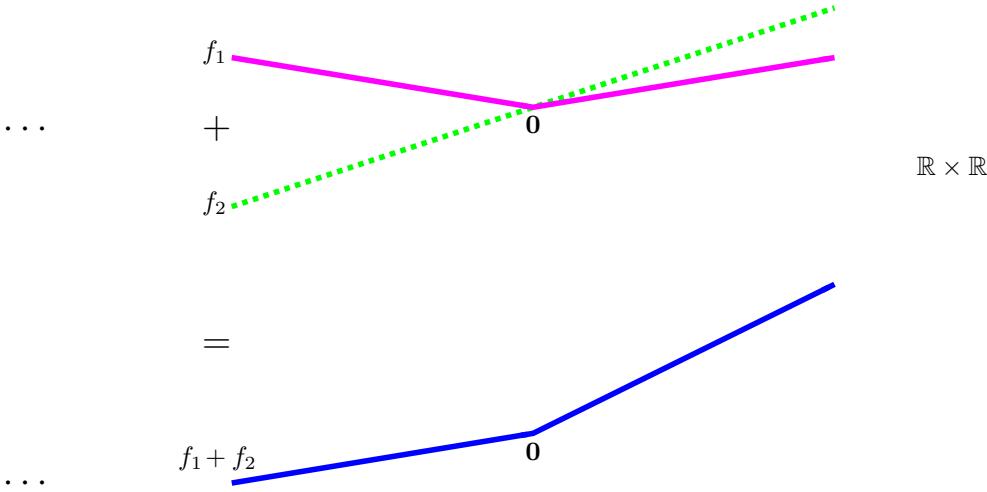


Figure 91: Nonnegatively weighted sum of convex functions is convex but can be unbounded below if any one function is. Unboundedness of the illustrated sum of real function  $f_1(x)=|x|$  with  $f_2(x)=ax$  depends upon slope of  $f_2$ . But  $g(x)=x^2+ax$  is never unbounded below, for any slope  $a$ , achieving its minimum at  $x=-a/2$ .

### 3.15 Salient properties of convex and quasiconvex functions

1. • A convex function is assumed continuous but not necessarily differentiable on the relative interior of its domain. [343, §10]
- A quasiconvex function is not necessarily a continuous function.
2. convex epigraph  $\Leftrightarrow$  convexity  $\Rightarrow$  quasiconvexity  $\Leftrightarrow$  convex sublevel sets.  
convex hypograph  $\Leftrightarrow$  concavity  $\Rightarrow$  quasiconcavity  $\Leftrightarrow$  convex superlevel sets.  
quasilinearity  $\Leftrightarrow$  convex level sets.
3. log-convex  $\Rightarrow$  convex  $\Rightarrow$  quasiconvex. 3.33  
concave  $\Rightarrow$  quasiconcave  $\Leftarrow$  log-concave  $\Leftarrow$  positive concave.
4. *Line Theorem* 3.13.0.0.1 translates identically to quasiconvexity (quasiconcavity). [127]
5. •  $g$  convex  $\Leftrightarrow -g$  concave.  
 $g$  quasiconvex  $\Leftrightarrow -g$  quasiconcave.  
 $g$  log-convex  $\Leftrightarrow 1/g$  log-concave.
- (*translation, homogeneity*) Function convexity, concavity, quasiconvexity, and quasiconcavity are invariant to offset and nonnegative scaling.
6. (*affine transformation of argument*) Composition  $g(h(X))$  of (quasi) convex (concave) function  $g$  with any affine function  $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times k}$  remains (quasi) convex (concave) in  $X \in \mathbb{R}^{m \times n}$ , where  $h(\mathbb{R}^{m \times n}) \cap \text{dom } g \neq \emptyset$ . [225, §B.2.1]

<sup>3.33</sup>Log-convex means: logarithm of function  $f$  is convex on  $\text{dom } f$ .

7. a. i. Nonnegatively weighted sum of convex (concave) functions remains convex (concave).<sup>3.34</sup> (Figure 91, §3.1.1.2.1)
- ii. Nonnegatively weighted max (min) of convex<sup>3.35</sup> (concave) functions remains convex (concave).
- iii. Pointwise supremum (infimum) of convex (concave) functions remains convex (concave). (Figure 78) [343, §5]
- iv.  $g$  convex nondecreasing monotonic **and**  $h$  convex  $\Rightarrow g(h)$  is convex.  
 $g$  concave nondecreasing monotonic **and**  $h$  concave  $\Rightarrow g(h)$  is concave. (§3.6.1.0.4)
- b. i. Sum of quasiconvex functions is not necessarily quasiconvex.
- ii. Nonnegatively weighted max (min) of quasiconvex (quasiconcave) functions remains quasiconvex (quasiconcave).
- iii. Pointwise supremum (infimum) of quasiconvex (quasiconcave) functions remains quasiconvex (quasiconcave).
- iv.  $g$  nondecreasing monotonic **and**  $h$  quasiconvex  $\Rightarrow g(h)$  is quasiconvex.  
 $g$  nondecreasing monotonic **and**  $h$  (quasi)concave  $\Rightarrow g(h)$  is quasiconcave.  
*[sic]* [74, thm.2.2.6] [378, thm.8.5]

#### 3.15.0.0.1 Exercise. Quasicomposition.

Fill in the blanks:

$g$  nonincreasing monotonic **and**  $h$  quasiconcave  $\Rightarrow g(h)$  is \_\_\_\_\_.

$g$  nonincreasing monotonic **and**  $h$  quasiconvex  $\Rightarrow g(h)$  is \_\_\_\_\_.




---

<sup>3.34</sup>Nonnegatively weighted nonzero sum of strictly convex (concave) functions remains strictly convex (concave).

<sup>3.35</sup>Supremum and maximum of convex functions are proven convex by intersection of epigraphs.



## Chapter 4

# Semidefinite programming

Prior to 1984, linear and nonlinear programming,<sup>4.1</sup> one a subset of the other, had evolved for the most part along unconnected paths, without even a common terminology. (The use of “programming” to mean “optimization” serves as a persistent reminder of these differences.)

—Forsgren, Gill, & Wright, 2002 [165]

Given some practical application of convex analysis, it may at first seem puzzling why a search for its solution ends abruptly with a formalized statement of the problem itself as a constrained optimization. The explanation is: typically we do not seek analytical solution because there are relatively few. (§3.5.3, §C) If a problem can be expressed in convex form, rather, then there exist computer programs providing efficient numerical global solution. [191] [446] [447] [445] [389] [373] The goal, then, becomes conversion of a given problem (perhaps a nonconvex or combinatorial problem statement) to an equivalent convex form or to an *alternation* of convex subproblems convergent to a solution of the original problem:

By the *fundamental theorem of Convex Optimization*, any locally optimal point (solution) of a convex problem is globally optimal. [65, §4.2.2] [342, §1] Given convex real objective function  $g$  and convex feasible set  $\mathcal{D} \subseteq \text{dom } g$ , which is the set of all variable values satisfying the problem constraints, we pose a generic convex optimization problem

$$\begin{aligned} &\underset{X}{\text{minimize}} && g(X) \\ &\text{subject to} && X \in \mathcal{D} \end{aligned} \tag{694}$$

where constraints are abstract here in membership of variable  $X$  to convex feasible set  $\mathcal{D}$ . Inequality constraint functions of a convex optimization problem are convex while equality constraint functions are conventionally affine, but not necessarily so. Affine equality constraint functions, as opposed to the superset of all convex equality constraint functions having convex level sets (§3.4.0.4), make convex optimization tractable.

Similarly, the problem

$$\begin{aligned} &\underset{X}{\text{maximize}} && g(X) \\ &\text{subject to} && X \in \mathcal{D} \end{aligned} \tag{695}$$

is called *convex* were  $g$  a real concave function and feasible set  $\mathcal{D}$  convex. As conversion to convex form is not always possible, there is much ongoing research to determine which problem classes have convex expression or relaxation. [35] [63] [172] [309] [384] [169]

---

<sup>4.1</sup> nascence of polynomial-time *interior-point methods* of solution [404] [443].  
Linear programming  $\subset$  (convex  $\cap$  nonlinear) programming.

## 4.1 Conic problem

*Still, we are surprised to see the relatively small number of submissions to semidefinite programming (SDP) solvers, as this is an area of significant current interest to the optimization community. We speculate that semidefinite programming is simply experiencing the fate of most new areas: Users have yet to understand how to pose their problems as semidefinite programs, and the lack of support for SDP solvers in popular modelling languages likely discourages submissions.*

— SIAM News, 2002. [133, p.9]

(confer p.127) Consider a *conic problem* (p) and its dual (d): [327, §3.3.1] [268, §2.1] [269]

$$(306) \quad \begin{array}{ll} \text{minimize}_{\substack{x}} & c^T x \\ \text{(p)} & \text{subject to } x \in \mathcal{K} \\ & Ax = b \end{array} \quad \begin{array}{ll} \text{maximize}_{\substack{y, s}} & b^T y \\ \text{(d)} & \text{subject to } s \in \mathcal{K}^* \\ & A^T y + s = c \end{array} \quad (696)$$

where  $\mathcal{K}$  is a closed convex cone,  $\mathcal{K}^*$  is its dual, matrix  $A$  is fixed, and the remaining quantities are vectors.

When  $\mathcal{K}$  is a polyhedral cone (§2.12.1), then each conic problem becomes a *linear program*; the selfdual nonnegative orthant providing the prototypical primal linear program and its dual. [103, §3-1]<sup>4.2</sup> More generally, each optimization problem is convex when  $\mathcal{K}$  is a closed convex cone. Solution to each convex problem is not necessarily unique; the optimal solution sets  $\{x^*\}$  and  $\{y^*, s^*\}$  are convex and may comprise more than a single point.

### 4.1.1 a semidefinite program

When  $\mathcal{K}$  is the selfdual cone of positive semidefinite matrices  $\mathbb{S}_+^n$  in the subspace of symmetric matrices  $\mathbb{S}^n$ , then each conic problem is called *semidefinite program* (SDP); [309, §6.4] primal problem (P) having matrix variable  $X \in \mathbb{S}^n$  while corresponding dual (D) has *slack variable*  $S \in \mathbb{S}^n$  and vector variable  $y = [y_i] \in \mathbb{R}^m$ : [11] [12, §2] [453, §1.3.8]

$$(P) \quad \begin{array}{ll} \text{minimize}_{\substack{X \in \mathbb{S}^n}} & \langle C, X \rangle \\ \text{subject to} & X \succeq 0 \\ & A \text{ svec } X = b \end{array} \quad \begin{array}{ll} \text{maximize}_{\substack{y \in \mathbb{R}^m, S \in \mathbb{S}^n}} & \langle b, y \rangle \\ \text{subject to} & S \succeq 0 \\ & \text{svec}^{-1}(A^T y) + S = C \end{array} \quad (D) \quad (697)$$

This is the *prototypical primal semidefinite program* and its dual, where matrix  $C \in \mathbb{S}^n$  and vector  $b \in \mathbb{R}^m$  are fixed as is

$$A \triangleq \begin{bmatrix} \text{svec}(A_1)^T \\ \vdots \\ \text{svec}(A_m)^T \end{bmatrix} \in \mathbb{R}^{m \times n(n+1)/2} \quad (698)$$

because  $\{A_i \in \mathbb{S}^n, i=1 \dots m\}$  is given. Thus

$$\begin{aligned} A \text{ svec } X &= \begin{bmatrix} \langle A_1, X \rangle \\ \vdots \\ \langle A_m, X \rangle \end{bmatrix} \\ \text{svec}^{-1}(A^T y) &= \sum_{i=1}^m y_i A_i \end{aligned} \quad (699)$$

---

<sup>4.2</sup>Dantzig explains reasoning behind a nonnegativity constraint: ... negative quantities of activities are not possible. ... a negative number of cases cannot be shipped.

The vector inner-product for matrices is defined in the Euclidean/Frobenius sense in the isomorphic vector space  $\mathbb{R}^{n(n+1)/2}$ ; *id est*,

$$\langle C, X \rangle \triangleq \text{tr}(C^T X) = \text{svec}(C)^T \text{svec } X \quad (38)$$

where  $\text{svec } X$  defined by (57) denotes symmetric vectorization.

*In a national planning problem of some size, one may easily run into several hundred variables and perhaps a hundred or more degrees of freedom. ... It should always be remembered that any mathematical method and particularly methods in linear programming must be judged with reference to the type of computing machinery available. Our outlook may perhaps be changed when we get used to the super modern, high capacity electronic computor that will be available here from the middle of next year.*

—Ragnar Frisch [167]

The Simplex method of solution for linear programming, invented by Dantzig in 1947 [103], is now integral to modern technology. The same cannot yet be said for semidefinite programming whose roots trace back to systems of positive semidefinite linear inequalities studied by Bellman & Fan in 1963 [32] [114] who provided saddle convergence criteria. Interior-point methods for numerical solution of linear programs can be traced back to the logarithmic barrier of Frisch in 1954 and Fiacco & McCormick in 1968 [160]. Karmarkar's polynomial-time interior-point method sparked a log-barrier renaissance in 1984, [306, §11] [443] [404] [309, p.3] but numerical performance of contemporary general-purpose semidefinite program solvers remains limited: Computational intensity for dense systems varies as  $O(m^2n)$  (*m constraints  $\ll n$  variables*) based on interior-point methods that produce solutions no more relatively accurate than 1E-8. There are no solvers capable of handling in excess of  $n=100,000$  variables without significant, sometimes crippling, loss of precision or time.<sup>4.3</sup> [36] [308, p.258] [73, p.3]

Nevertheless, semidefinite programming has recently emerged to prominence because it admits a new class of problem previously unsolvable by convex optimization techniques, [63] and because it theoretically subsumes other convex techniques: (Figure 92) linear programming and *quadratic programming* and *second-order cone programming*.<sup>4.4</sup> Determination of the Riemann mapping function from complex analysis [318] [30, §8, §13], for example, can be posed as a semidefinite program.

### 4.1.2 Maximal complementarity

It has been shown [453, §2.5.3] that contemporary interior-point methods [444] [321] [309] [12] [65, §11] [165] (developed *circa* 1990 [172] for numerical solution of semidefinite programs) can converge to a solution of *maximal complementarity*; [200, §5] [452] [283] [179] not a vertex solution but a solution of highest cardinality or rank among all optimal solutions.<sup>4.5</sup>

This phenomenon can be explained by recognizing that interior-point methods generally find solutions relatively interior to a feasible set by design.<sup>4.6</sup> [7, p.3] Log barriers are designed to fail numerically at the feasible set boundary. So low-rank solutions, all

<sup>4.3</sup>Heuristics are not ruled out by SIOPT; indeed I would suspect that most successful methods have (appropriately described) heuristics under the hood - my codes certainly do. ... Of course, there are still questions relating to high-accuracy and speed, but for many applications a few digits of accuracy suffices and overnight runs for non-real-time delivery is acceptable.

—Nicholas I. M. Gould, [Stanford alumnus](#), SIOPT Editor in Chief

<sup>4.4</sup>Second-order cone programming was born in the 1990s; it is not posable as a quadratic program. [278]

<sup>4.5</sup>This characteristic might be regarded as a disadvantage to interior-point methods of numerical solution, but this behavior is not certain and depends on solver implementation.

<sup>4.6</sup>Simplex methods, in contrast, find vertex solutions. [103, p.158] [16, p.2]

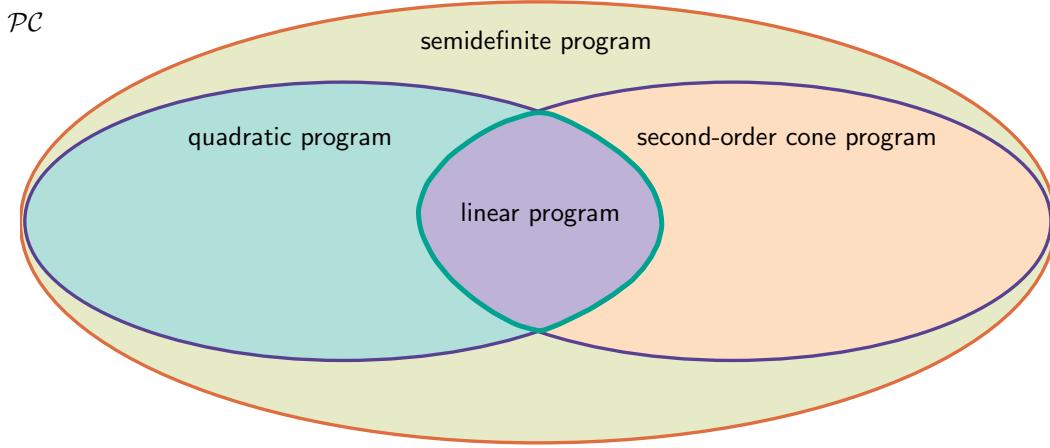


Figure 92: Venn diagram of programming hierarchy. Convex program  $\mathcal{PC}$  represents the broadest class of optimization problem for which efficient solution exists theoretically. ( $\setminus \mathcal{PC}$  comprises those for which convex equivalents have not yet been found.) Semidefinite program, a subset of  $\mathcal{PC}$ , subsumes other convex program classes excepting geometric program but includes quadratically constrained program. Second-order cone program and quadratic program each subsume linear program.

on the boundary, are rendered more difficult to find as numerical error becomes more prevalent there.

#### 4.1.2.1 Reduced-rank solution

A simple rank reduction algorithm, for construction of a primal optimal solution  $X^*$  to (697P) satisfying an upper bound on rank governed by Proposition 2.9.3.0.1, is presented in §4.3. That proposition asserts existence of feasible solutions with an upper bound on their rank; [27, §II.13.1] specifically, it asserts an extreme point (§2.6.0.0.1) of *primal feasible set*  $\mathcal{A} \cap \mathbb{S}_+^n$  satisfies upper bound

$$\text{rank } X \leq \left\lfloor \frac{\sqrt{8m+1}-1}{2} \right\rfloor \quad (275)$$

where, given  $A \in \mathbb{R}^{m \times n(n+1)/2}$  (698) and  $b \in \mathbb{R}^m$ ,

$$\mathcal{A} \triangleq \{X \in \mathbb{S}^n \mid A \text{ svec } X = b\} \quad (2273)$$

is the affine subset from primal problem (697P).

#### 4.1.2.2 Coexistence of low- and high-rank solutions; analogy

That low-rank and high-rank optimal solutions  $\{X^*\}$  of (697P) coexist may be grasped with the following analogy: We compare a proper polyhedral cone  $\mathcal{S}_+^3$  in  $\mathbb{R}^3$  (illustrated in Figure 93) to the positive semidefinite cone  $\mathbb{S}_+^3$  in isometrically isomorphic  $\mathbb{R}^6$ , difficult to visualize. The analogy is good:

- $\text{intr } \mathbb{S}_+^3$  is constituted by rank-3 matrices.
- $\text{intr } \mathcal{S}_+^3$  has three dimensions.

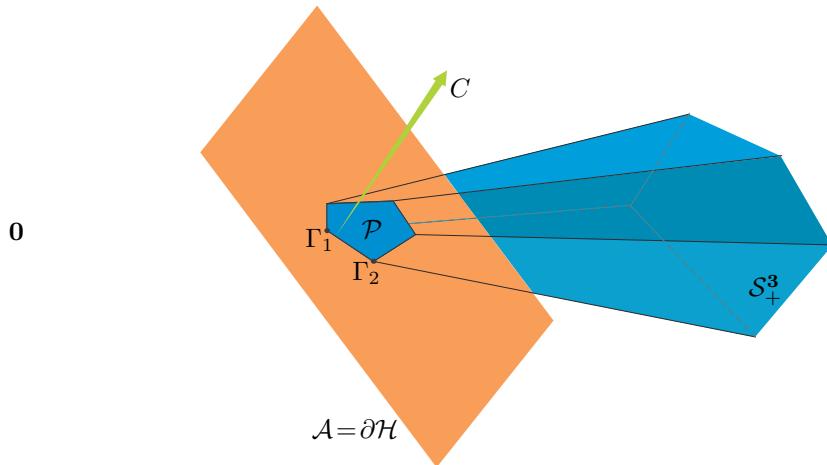


Figure 93: Visualizing positive semidefinite cone in high dimension: Proper polyhedral cone  $\mathcal{S}_+^3 \subset \mathbb{R}^3$  representing positive semidefinite cone  $\mathbb{S}_+^3 \subset \mathbb{S}^3$ ; analogizing its intersection  $\mathbb{S}_+^3 \cap \partial\mathcal{H}$  with hyperplane. Number of facets is arbitrary (an analogy not inspired by eigenvalue decomposition). The rank-0 positive semidefinite matrix corresponds to origin in  $\mathbb{R}^3$ , rank-1 positive semidefinite matrices correspond to edges of polyhedral cone, rank-2 to facet relative interiors, and rank-3 to polyhedral cone interior. Vertices  $\Gamma_1$  and  $\Gamma_2$  are extreme points of polyhedron  $\mathcal{P} = \partial\mathcal{H} \cap \mathcal{S}_+^3$ , and extreme directions of  $\mathcal{S}_+^3$ . A given vector  $C$  is normal to another hyperplane (not illustrated but independent w.r.t  $\partial\mathcal{H}$ ) containing line segment  $\overline{\Gamma_1\Gamma_2}$  minimizing real linear function  $\langle C, X \rangle$  on  $\mathcal{P}$ . (confer Figure 29, Figure 33)

- boundary  $\partial \mathbb{S}_+^3$  contains rank-0, rank-1, and rank-2 matrices.  
boundary  $\partial \mathcal{S}_+^3$  contains 0-, 1-, and 2-dimensional faces.
- the only rank-0 matrix resides in the vertex at the origin.
- Rank-1 matrices are in one-to-one correspondence with extreme directions of  $\mathbb{S}_+^3$  and  $\mathcal{S}_+^3$ . The set of all rank-1 symmetric matrices in this dimension

$$\{G \in \mathbb{S}_+^3 \mid \text{rank } G = 1\} \quad (700)$$

is not a connected set.

- Rank of a sum of members  $F+G$  in Lemma 2.9.2.9.1 and location of a difference  $F-G$  in §2.9.2.12.1 similarly hold for  $\mathbb{S}_+^3$  and  $\mathcal{S}_+^3$ .
- Euclidean distance from any particular rank-3 positive semidefinite matrix (in the cone interior) to the closest rank-2 positive semidefinite matrix (on the boundary) is generally less than the distance to the closest rank-1 positive semidefinite matrix. (§7.1.2)
- distance from any point in  $\partial \mathbb{S}_+^3$  to  $\text{intr } \mathbb{S}_+^3$  is infinitesimal (§2.1.7.1.1).  
distance from any point in  $\partial \mathcal{S}_+^3$  to  $\text{intr } \mathcal{S}_+^3$  is infinitesimal.
- faces of  $\mathbb{S}_+^3$  correspond to faces of  $\mathcal{S}_+^3$  (confer Table 2.9.2.3.1):

	$k$	$\dim \mathcal{F}(\mathcal{S}_+^3)$	$\dim \mathcal{F}(\mathbb{S}_+^3)$	$\dim \mathcal{F}(\mathbb{S}_+^3 \ni \text{rank-}k \text{ matrix})$
boundary	0	0	0	0
	1	1	1	1
	2	2	3	3
interior	3	3	6	6

Integer  $k$  indexes  $k$ -dimensional faces  $\mathcal{F}$  of  $\mathcal{S}_+^3$ . Positive semidefinite cone  $\mathbb{S}_+^3$  has four kinds of faces, including cone itself ( $k=3$ , boundary + interior), whose dimensions in isometrically isomorphic  $\mathbb{R}^6$  are listed under  $\dim \mathcal{F}(\mathbb{S}_+^3)$ . Smallest face  $\mathcal{F}(\mathbb{S}_+^3 \ni \text{rank-}k \text{ matrix})$  that contains a rank- $k$  positive semidefinite matrix has dimension  $k(k+1)/2$  by (226).

- For  $\mathcal{A}$  equal to intersection of  $m$  hyperplanes having linearly independent normals, and for  $X \in \mathcal{S}_+^3 \cap \mathcal{A}$ , we have  $\text{rank } X \leq m$ ; the analogue to (275).

**Proof.** With reference to Figure 93: Assume one ( $m=1$ ) hyperplane  $\mathcal{A}=\partial \mathcal{H}$  intersects the polyhedral cone. Every intersecting plane contains at least one matrix having rank less than or equal to 1; *id est*, from all  $X \in \partial \mathcal{H} \cap \mathcal{S}_+^3$  there exists an  $X$  such that  $\text{rank } X \leq 1$ . Rank 1 is therefore an upper bound in this case.

Now visualize intersection of the polyhedral cone with two ( $m=2$ ) hyperplanes having linearly independent normals. The hyperplane intersection  $\mathcal{A}$  makes a line. Every intersecting line contains at least one matrix having rank less than or equal to 2, providing an upper bound. In other words, there exists a positive semidefinite matrix  $X$  belonging to any line intersecting the polyhedral cone such that  $\text{rank } X \leq 2$ .

In the case of three independent intersecting hyperplanes ( $m=3$ ), the hyperplane intersection  $\mathcal{A}$  makes a point that can reside anywhere in the polyhedral cone. The upper bound on a point in  $\mathcal{S}_+^3$  is also the greatest upper bound:  $\text{rank } X \leq 3$ . ♦

#### 4.1.2.2.1 Example. Optimization over $\mathcal{A} \cap \mathcal{S}_+^3$ .

Consider minimization of the real linear function  $\langle C, X \rangle$  over

$$\mathcal{P} \triangleq \mathcal{A} \cap \mathcal{S}_+^3 \quad (701)$$

a polyhedral feasible set;

$$\begin{aligned} f_0^* &\triangleq \underset{X}{\text{minimize}} \quad \langle C, X \rangle \\ &\text{subject to} \quad X \in \mathcal{A} \cap \mathcal{S}_+^3 \end{aligned} \quad (702)$$

As illustrated for particular vector  $C$  and hyperplane  $\mathcal{A} = \partial\mathcal{H}$  in Figure 93, this linear function is minimized on any  $X$  belonging to the face of  $\mathcal{P}$  containing extreme points  $\{\Gamma_1, \Gamma_2\}$  and all the rank-2 matrices in between; *id est*, on any  $X$  belonging to the face of  $\mathcal{P}$

$$\mathcal{F}(\mathcal{P}) = \{X \mid \langle C, X \rangle = f_0^*\} \cap \mathcal{A} \cap \mathcal{S}_+^3 \quad (703)$$

exposed by the hyperplane  $\{X \mid \langle C, X \rangle = f_0^*\}$ . In other words, the set of all optimal points  $X^*$  is a face of  $\mathcal{P}$

$$\{X^*\} = \mathcal{F}(\mathcal{P}) = \overline{\Gamma_1 \Gamma_2} \quad (704)$$

comprising rank-1 and rank-2 positive semidefinite matrices. Rank 1 is the upper bound on existence in the feasible set  $\mathcal{P}$  for this case  $m=1$  hyperplane constituting  $\mathcal{A}$ . The rank-1 matrices  $\Gamma_1$  and  $\Gamma_2$  in face  $\mathcal{F}(\mathcal{P})$  are extreme points of that face and (by transitivity (§2.6.1.2)) extreme points of the intersection  $\mathcal{P}$  as well. As predicted by analogy to Barvinok's Proposition 2.9.3.0.1, the upper bound on rank of  $X$  existent in the feasible set  $\mathcal{P}$  is satisfied by an extreme point. The upper bound on rank of an optimal solution  $X^*$  existent in  $\mathcal{F}(\mathcal{P})$  is thereby also satisfied by an extreme point of  $\mathcal{P}$  precisely because  $\{X^*\}$  constitutes  $\mathcal{F}(\mathcal{P})$ ; <sup>4.7</sup> in particular,

$$\{X^* \in \mathcal{P} \mid \text{rank } X^* \leq 1\} = \{\Gamma_1, \Gamma_2\} \subseteq \mathcal{F}(\mathcal{P}) \quad (705)$$

As all linear functions on a polyhedron are minimized on a face, [103] [282] [305] [312] by analogy we so demonstrate coexistence of optimal solutions  $X^*$  of (697P) having assorted rank.  $\square$

#### 4.1.2.3 Previous work

Barvinok showed, [25, §2.2] when given a positive definite matrix  $C$  and an arbitrarily small neighborhood of  $C$  comprising positive definite matrices, there exists a matrix  $\tilde{C}$  from that neighborhood such that optimal solution  $X^*$  to (697P) (substituting  $\tilde{C}$ ) is an extreme point of  $\mathcal{A} \cap \mathbb{S}_+^n$  and satisfies upper bound (275).<sup>4.8</sup> Given arbitrary positive definite  $C$ , this means nothing inherently guarantees that an optimal solution  $X^*$  to problem (697P) satisfies (275); certainly nothing given any symmetric matrix  $C$ , as the problem is posed. This can be proved by example:

<sup>4.7</sup> and every face contains a subset of the extreme points of  $\mathcal{P}$  by the *extreme existence theorem* (§2.6.0.0.2). This means: because the affine subset  $\mathcal{A}$  and hyperplane  $\{X \mid \langle C, X \rangle = f_0^*\}$  must intersect a whole face of  $\mathcal{P}$ , calculation of an upper bound on rank of  $X^*$  ignores counting the hyperplane when determining  $m$  in (275).

<sup>4.8</sup> Further, the set of all such  $\tilde{C}$  in that neighborhood is open and dense.

#### 4.1.2.3.1 Example. (Ye) Maximal Complementarity.

Assume dimension  $n$  to be an even positive number. Then the particular instance of problem (697P),

$$\begin{array}{ll} \text{minimize}_{X \in \mathbb{S}^n} & \left\langle \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & 2I \end{bmatrix}, X \right\rangle \\ \text{subject to} & X \succeq 0 \\ & \langle I, X \rangle = n \end{array} \quad (706)$$

has optimal solution

$$X^* = \begin{bmatrix} 2I & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{S}^n \quad (707)$$

with an equal number of twos and zeros along the main diagonal. Indeed, optimal solution (707) is a terminal solution along the *central path* taken by the interior-point method as implemented in [453, §2.5.3]; it is also a solution of highest rank among all optimal solutions to (706). Clearly, rank of this primal optimal solution exceeds by far a rank-1 solution predicted by upper bound (275).  $\square$

#### 4.1.2.4 Later developments

This rational example (706) indicates the need for a more generally applicable and simple algorithm to identify an optimal solution  $X^*$  satisfying Barvinok's Proposition 2.9.3.0.1. We will review such an algorithm in §4.3, but first we provide more background.

## 4.2 Framework

### 4.2.1 Feasible sets

Denote by  $\mathcal{D}$  and  $\mathcal{D}^*$  the convex sets of primal and dual points respectively satisfying the primal and dual constraints in (697), each assumed nonempty;

$$\begin{aligned} \mathcal{D} &= \left\{ X \in \mathbb{S}_+^n \mid \begin{bmatrix} \langle A_1, X \rangle \\ \vdots \\ \langle A_m, X \rangle \end{bmatrix} = b \right\} = \mathcal{A} \cap \mathbb{S}_+^n \\ \mathcal{D}^* &= \left\{ S \in \mathbb{S}_+^n, y = [y_i] \in \mathbb{R}^m \mid \sum_{i=1}^m y_i A_i + S = C \right\} \end{aligned} \quad (708)$$

These are the *primal feasible set* and *dual feasible set*. Geometrically, primal feasible  $\mathcal{A} \cap \mathbb{S}_+^n$  represents an intersection of the positive semidefinite cone  $\mathbb{S}_+^n$  with an affine subset  $\mathcal{A}$  of the subspace of symmetric matrices  $\mathbb{S}^n$  in isometrically isomorphic  $\mathbb{R}^{n(n+1)/2}$ .  $\mathcal{A}$  has dimension  $n(n+1)/2 - m$  when the vectorized  $A_i$  are linearly independent. Dual feasible set  $\mathcal{D}^*$  is a Cartesian product of the positive semidefinite cone with its inverse image (§2.1.9.0.1) under affine transformation<sup>4.9</sup>  $C - \sum y_i A_i$ . Both feasible sets are convex, and the objective functions are linear on a Euclidean vector space. Hence, (697P) and (697D) are convex optimization problems.

---

<sup>4.9</sup>Inequality  $C - \sum y_i A_i \succeq 0$  follows directly from (697D) (§2.9.0.1.1) and is known as a *linear matrix inequality*. (§2.13.6.1.1) Because  $\sum y_i A_i \preceq C$ , matrix  $S$  is known as a *slack variable* (a term borrowed from linear programming [103]) since its inclusion raises this inequality to equality.

#### 4.2.1.1 $\mathcal{A} \cap \mathbb{S}_+^n$ emptiness determination via Farkas' lemma

**4.2.1.1.1 Lemma.** *Semidefinite Farkas' lemma.* (confer §4.2.1.1.2)

Given affine subset  $\mathcal{A} = \{X \in \mathbb{S}^n \mid \langle A_i, X \rangle = b_i, i=1 \dots m\}$  (2273), vector  $b = [b_i] \in \mathbb{R}^m$ , and set  $\{A_i \in \mathbb{S}^n, i=1 \dots m\}$  such that  $\{A \text{ svec } X \mid X \succeq 0\}$  (386) is closed, then primal feasible set  $\mathcal{A} \cap \mathbb{S}_+^n$  is nonempty if and only if  $y^T b \geq 0$  holds for each and every vector  $y = [y_i] \in \mathbb{R}^m$  such that  $\sum_{i=1}^m y_i A_i \succeq 0$ .

Equivalently, primal feasible set  $\mathcal{A} \cap \mathbb{S}_+^n$  is nonempty if and only if  $y^T b \geq 0$  holds for each and every vector  $\|y\|=1$  such that  $\sum_{i=1}^m y_i A_i \succeq 0$ .  $\diamond$

*Semidefinite Farkas' lemma* provides necessary and sufficient conditions for a set of hyperplanes to have nonempty intersection  $\mathcal{A} \cap \mathbb{S}_+^n$  with the positive semidefinite cone. Given

$$A = \begin{bmatrix} \text{svec}(A_1)^T \\ \vdots \\ \text{svec}(A_m)^T \end{bmatrix} \in \mathbb{R}^{m \times n(n+1)/2} \quad (698)$$

*semidefinite Farkas' lemma* assumes that a convex cone

$$\mathcal{K} = \{A \text{ svec } X \mid X \succeq 0\} \quad (386)$$

is closed per membership relation (323) from which the lemma springs: [260, §I]  $\mathcal{K}$  closure is attained when matrix  $A$  satisfies the *cone closedness invariance corollary* (p.143). Given closed convex cone  $\mathcal{K}$  and its dual from Example 2.13.6.1.1

$$\mathcal{K}^* = \{y \mid \sum_{j=1}^m y_j A_j \succeq 0\} \quad (393)$$

then we can apply membership relation

$$b \in \mathcal{K} \Leftrightarrow \langle y, b \rangle \geq 0 \quad \forall y \in \mathcal{K}^* \quad (323)$$

to obtain the lemma

$$b \in \mathcal{K} \Leftrightarrow \exists X \succeq 0 \ni A \text{ svec } X = b \Leftrightarrow \mathcal{A} \cap \mathbb{S}_+^n \neq \emptyset \quad (709)$$

$$b \in \mathcal{K} \Leftrightarrow \langle y, b \rangle \geq 0 \quad \forall y \in \mathcal{K}^* \Leftrightarrow \mathcal{A} \cap \mathbb{S}_+^n \neq \emptyset \quad (710)$$

The final equivalence synopsizes *semidefinite Farkas' lemma*.

While the lemma is correct as stated, a positive definite version is required for semidefinite programming [453, §1.3.8] because existence of a feasible solution in the cone interior  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  is required by *Slater's condition*<sup>4.10</sup> to achieve 0 duality gap (optimal primal–dual objective difference, §4.2.3, Figure 64). Geometrically, a positive definite lemma is required to insure that a point of intersection closest to the origin is not at infinity; e.g., Figure 48. Then given  $A \in \mathbb{R}^{m \times n(n+1)/2}$  having rank  $m$ , we wish to detect existence of nonempty primal feasible set interior to the PSD cone;<sup>4.11</sup> (389)

$$b \in \text{intr } \mathcal{K} \Leftrightarrow \langle y, b \rangle > 0 \quad \forall y \in \mathcal{K}^*, \quad y \neq \mathbf{0} \Leftrightarrow \mathcal{A} \cap \text{intr } \mathbb{S}_+^n \neq \emptyset \quad (711)$$

*Positive definite Farkas' lemma* is made from proper cones,  $\mathcal{K}$  (386) and  $\mathcal{K}^*$  (393), and membership relation (329) for which  $\mathcal{K}$  closedness is unnecessary:

<sup>4.10</sup>Slater's sufficient constraint qualification is satisfied whenever any primal or dual *strictly feasible solution* exists; *id est*, any point satisfying the respective affine constraints and relatively interior to the convex cone. [366, §6.6] [42, p.325] If the cone were polyhedral, then Slater's constraint qualification is satisfied when any feasible solution exists (relatively interior to the cone or on its relative boundary). [65, §5.2.3]

<sup>4.11</sup>Detection of  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n \neq \emptyset$  by examining  $\text{intr } \mathcal{K}$  instead is a trick need not be lost.

**4.2.1.1.2 Lemma.** *Positive definite Farkas' lemma.* (confer §4.2.1.1.1)  
Given l.i. set  $\{A_i \in \mathbb{S}^n, i=1 \dots m\}$  and vector  $b = [b_i] \in \mathbb{R}^m$ , make affine set

$$\mathcal{A} = \{X \in \mathbb{S}^n \mid \langle A_i, X \rangle = b_i, i=1 \dots m\} \quad (2273)$$

Primal feasible cone interior  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  is nonempty if and only if  $y^T b > 0$  holds for each and every vector  $y = [y_i] \neq \mathbf{0}$  such that  $\sum_{i=1}^m y_i A_i \succeq 0$ .

Equivalently, primal feasible cone interior  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  is nonempty if and only if  $y^T b > 0$  holds for each and every vector  $\|y\|=1 \Rightarrow \sum_{i=1}^m y_i A_i \succeq 0$ .  $\diamond$

#### 4.2.1.1.3 Example. “New” Farkas’ lemma.

Lasserre [260, §III] presented an example in 1995, originally offered by Ben-Israel in 1969 [33, p.378], to support closedness in *semidefinite Farkas’ Lemma 4.2.1.1.1*:

$$A \triangleq \begin{bmatrix} \text{svec}(A_1)^T \\ \text{svec}(A_2)^T \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (712)$$

Intersection  $\mathcal{A} \cap \mathbb{S}_+^n$  is practically empty because the solution set

$$\{X \succeq 0 \mid A \text{ svec } X = b\} = \left\{ \begin{bmatrix} \alpha & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & 0 \end{bmatrix} \succeq 0 \mid \alpha \in \mathbb{R} \right\} \quad (713)$$

is positive semidefinite only asymptotically ( $\alpha \rightarrow \infty$ ). Yet  $\sum_{i=1}^m y_i A_i \succeq 0 \Rightarrow y^T b \geq 0$  the dual system erroneously indicates nonempty intersection because  $\mathcal{K}$  (386) violates a closedness condition of the lemma; *videlicet*, for  $\|y\|=1$

$$y_1 \begin{bmatrix} 0 & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & 0 \end{bmatrix} + y_2 \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \succeq 0 \Leftrightarrow y = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \Rightarrow y^T b = 0 \quad (714)$$

On the other hand, *positive definite Farkas’ Lemma 4.2.1.1.2* certifies that  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  is empty; what we need to know for semidefinite programming.

Lasserre suggested addition of another condition to *semidefinite Farkas’ lemma* (§4.2.1.1.1) to make a new lemma having no closedness condition. But *positive definite Farkas’ lemma* (§4.2.1.1.2) is simpler and obviates the additional condition proposed.  $\square$

#### 4.2.1.2 Theorem of the alternative for semidefinite programming

Because these Farkas’ lemmas follow from membership relations, we may construct alternative systems from them. Applying the method of §2.13.2.1.1, then from *positive definite Farkas’ lemma* we get

$$\begin{aligned} \mathcal{A} \cap \text{intr } \mathbb{S}_+^n &\neq \emptyset \\ \text{or in the alternative} \\ y^T b \leq 0, \quad \sum_{i=1}^m y_i A_i &\succeq 0, \quad y \neq \mathbf{0} \end{aligned} \quad (715)$$

Any single vector  $y$  satisfying the alternative certifies  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  is empty. Such a vector can be found as a solution to another semidefinite program: for linearly independent

(vectorized) set  $\{A_i \in \mathbb{S}^n, i=1 \dots m\}$

$$\begin{aligned} & \underset{y}{\text{minimize}} && y^T b \\ & \text{subject to} && \sum_{i=1}^m y_i A_i \succeq 0 \\ & && \|y\|^2 \leq 1 \end{aligned} \tag{716}$$

If an optimal vector  $y^* \neq \mathbf{0}$  can be found such that  $y^{*\top} b \leq 0$ , then primal feasible cone interior  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  is empty.

#### 4.2.1.3 Boundary-membership criterion

(confer (710)(711)) From boundary-membership relation (333), for proper cones  $\mathcal{K}$  (386) and  $\mathcal{K}^*$  (393) of linear matrix inequality,

$$b \in \partial \mathcal{K} \Leftrightarrow \exists y \neq \mathbf{0} \ni \langle y, b \rangle = 0, \quad y \in \mathcal{K}^*, \quad b \in \mathcal{K} \Leftrightarrow \partial \mathbb{S}_+^n \cap \mathcal{A} \cap \mathbb{S}_+^n \neq \emptyset \tag{717}$$

Whether vector  $b \in \partial \mathcal{K}$  belongs to cone  $\mathcal{K}$  boundary, that is a determination we can indeed make; one that is certainly expressible as a feasibility problem: Given linearly independent set<sup>4.12</sup>  $\{A_i \in \mathbb{S}^n, i=1 \dots m\}$ , for  $b \in \mathcal{K}$  (709)

$$\begin{aligned} & \text{find} && y \neq \mathbf{0} \\ & \text{subject to} && y^T b = 0 \\ & && \sum_{i=1}^m y_i A_i \succeq 0 \end{aligned} \tag{718}$$

Any such nonzero solution  $y$  certifies that affine subset  $\mathcal{A}$  (2273) intersects the positive semidefinite cone  $\mathbb{S}_+^n$  only on its boundary; in other words, nonempty feasible set  $\mathcal{A} \cap \mathbb{S}_+^n$  belongs to the positive semidefinite cone boundary  $\partial \mathbb{S}_+^n$ .

#### 4.2.2 Duals

The dual objective function from (697D) evaluated at any feasible solution represents a lower bound on the primal optimal objective value from (697P). We can see this by direct substitution: Assume the feasible sets  $\mathcal{A} \cap \mathbb{S}_+^n$  and  $\mathcal{D}^*$  are nonempty. Then it is always true:

$$\begin{aligned} & \langle C, X \rangle \geq \langle b, y \rangle \\ & \left\langle \sum_i y_i A_i + S, X \right\rangle \geq [\langle A_1, X \rangle \cdots \langle A_m, X \rangle] y \\ & \langle S, X \rangle \geq 0 \end{aligned} \tag{719}$$

The converse also follows because

$$X \succeq 0, \quad S \succeq 0 \Rightarrow \langle S, X \rangle \geq 0 \tag{1639}$$

Optimal value of the dual objective thus represents the greatest lower bound on the primal. This fact is known as *weak duality* for semidefinite programming, [453, §1.3.8] and can be used to detect convergence in any primal/dual numerical method of solution.

---

<sup>4.12</sup>From the results of Example 2.13.6.1.1, vector  $b$  on the boundary of  $\mathcal{K}$  cannot be detected simply by looking for 0 eigenvalues in matrix  $X$ . We do not consider a thin-or-square matrix  $A$  because then feasible set  $\mathcal{A} \cap \mathbb{S}_+^n$  is at most a single point.

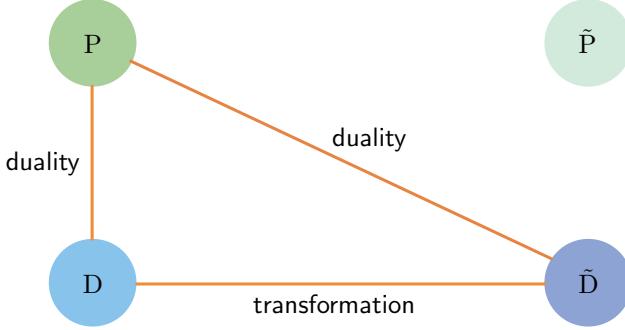


Figure 94: Connectivity indicates paths between particular primal and dual problems from Exercise 4.2.2.1.1. More generally, any path between primal problems  $P$  (and equivalent  $\tilde{P}$ ) and dual  $D$  (and equivalent  $\tilde{D}$ ) is possible: implying, any given path is not necessarily circuital; dual of a dual problem is not necessarily stated in precisely same manner as corresponding primal convex problem, in other words, although its solution set is equivalent to within some transformation.

#### 4.2.2.1 Dual problem statement is not unique

Even subtle but equivalent restatements of a primal convex problem can lead to vastly different statements of a corresponding dual problem. This phenomenon is of interest because a particular instantiation of dual problem might be easier to solve numerically or it might take one of few forms for which analytical solution is known.

Here is a canonical restatement of prototypical dual semidefinite program (697D), for example, equivalent by (198):

$$(D) \quad \begin{array}{ll} \underset{y \in \mathbb{R}^m, S \in \mathbb{S}^n}{\text{maximize}} & \langle b, y \rangle \\ \text{subject to} & S \succeq 0 \\ & \text{svec}^{-1}(A^T y) + S = C \end{array} \quad \equiv \quad \begin{array}{ll} \underset{y \in \mathbb{R}^m}{\text{maximize}} & \langle b, y \rangle \\ \text{subject to} & \text{svec}^{-1}(A^T y) \preceq C \end{array} \quad (697\tilde{D})$$

Dual feasible cone interior in  $\text{intr } \mathbb{S}_+^n$  (708) (699) thereby corresponds with canonical dual ( $\tilde{D}$ ) feasible interior

$$\text{rel intr } \tilde{D}^* \triangleq \left\{ y \in \mathbb{R}^m \mid \sum_{i=1}^m y_i A_i \prec C \right\} \quad (720)$$

##### 4.2.2.1.1 Exercise. Prototypical primal semidefinite program.

Derive prototypical primal (697P) from its canonical dual (697 $\tilde{D}$ ); *id est*, demonstrate that particular connectivity in Figure 94. ▼

#### 4.2.3 Optimality conditions

When primal feasible cone interior  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  exists in  $\mathbb{S}^n$  or when canonical dual feasible interior  $\text{rel intr } \tilde{D}^*$  exists in  $\mathbb{R}^m$ , then these two problems (697P) (697D) become strong duals by Slater's sufficient condition (p.229). In other words, the primal optimal objective value becomes equal to the dual optimal objective value: there is no duality gap (Figure 64) and so determination of convergence is facilitated; *id est*, if  $\exists X \in \mathcal{A} \cap \text{intr } \mathbb{S}_+^n$

or  $\exists y \in \text{rel intr } \tilde{\mathcal{D}}^*$  then

$$\begin{aligned} \langle C, X^* \rangle &= \langle b, y^* \rangle \\ \left\langle \sum_i y_i^* A_i + S^*, X^* \right\rangle &= [\langle A_1, X^* \rangle \cdots \langle A_m, X^* \rangle] y^* \\ \langle S^*, X^* \rangle &= 0 \end{aligned} \quad (721)$$

where  $S^*, y^*$  denote a dual optimal solution.<sup>4.13</sup> We summarize this:

**4.2.3.0.1 Corollary.** *Optimality and strong duality.* [400, §3.1] [453, §1.3.8]  
For semidefinite programs (697P) and (697D), assume primal and dual feasible sets  $\mathcal{A} \cap \mathbb{S}_+^n \subset \mathbb{S}^n$  and  $\mathcal{D}^* \subset \mathbb{S}^n \times \mathbb{R}^m$  (708) are nonempty. Then

- $X^*$  is optimal for (697P)
- $S^*, y^*$  are optimal for (697D)
- duality gap  $\langle C, X^* \rangle - \langle b, y^* \rangle$  is 0

if and only if

- i)  $\exists X \in \mathcal{A} \cap \text{intr } \mathbb{S}_+^n$  or  $\exists y \in \text{rel intr } \tilde{\mathcal{D}}^*$   
and
- ii)  $\langle S^*, X^* \rangle = 0$

◇

For symmetric positive semidefinite matrices, requirement ii is equivalent to the *complementarity*

$$\langle S^*, X^* \rangle = 0 \Leftrightarrow S^* X^* = X^* S^* = \mathbf{0} \quad (1748)$$

Commutativity of diagonalizable matrices is necessary and sufficient [228, §1.3.12] for these two optimal symmetric matrices to be simultaneously diagonalizable. Therefore

$$\text{rank } X^* + \text{rank } S^* \leq n \quad (722)$$

**Proof.** The product of symmetric optimal matrices  $X^*, S^* \in \mathbb{S}^n$  must itself be symmetric because of commutativity. (1628) The symmetric product has diagonalization [12, cor.2.11]

$$S^* X^* = X^* S^* = Q \Lambda_{S^*} \Lambda_{X^*} Q^T = \mathbf{0} \Leftrightarrow \Lambda_{X^*} \Lambda_{S^*} = \mathbf{0} \quad (723)$$

where  $Q$  is an orthogonal matrix. Product of the nonnegative diagonal  $\Lambda$  matrices can be  $\mathbf{0}$  if their main diagonal zeros are complementary or coincide. Due only to symmetry,  $\text{rank } X^* = \text{rank } \Lambda_{X^*}$  and  $\text{rank } S^* = \text{rank } \Lambda_{S^*}$  for these optimal primal and dual solutions. (1614) So total number of nonzero diagonal entries, from both  $\Lambda$ , cannot exceed  $n$  because of the complementarity. ♦

When equality is attained in (722)

$$\text{rank } X^* + \text{rank } S^* = n \quad (724)$$

there are no coinciding main diagonal zeros in  $\Lambda_{X^*} \Lambda_{S^*}$ , and so we have what is called *strict complementarity*.<sup>4.14</sup> Logically it follows that a necessary and sufficient condition for strict complementarity of an optimal primal and dual solution is

$$X^* + S^* \succ 0 \quad (725)$$

<sup>4.13</sup> Optimality condition  $\langle S^*, X^* \rangle = 0$  is called a *complementary slackness condition*, in keeping with LP tradition [103], that forbids dual inequalities in (697) to simultaneously hold strictly. [342, §4]

<sup>4.14</sup> distinct from maximal complementarity (§4.1.2).

#### 4.2.3.1 solving primal problem via dual

The beauty of Corollary 4.2.3.0.1 is its conjugacy; *id est*, one can solve either the primal or dual problem in (697) and then find a solution to the other via the optimality conditions. When a dual optimal solution is known, for example, a primal optimal solution is any primal feasible solution in hyperplane  $\{X \mid \langle S^*, X \rangle = 0\}$ .

**4.2.3.1.1 Example.** *Minimal cardinality Boolean.* [102] [35, §4.3.4] [384] (confer Example 4.6.1.5.1) Consider finding a *minimal cardinality* Boolean solution  $x$  to the classic linear algebra problem  $Ax = b$  given noiseless data  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ ;

$$\begin{aligned} & \underset{x}{\text{minimize}} \quad \|x\|_0 \\ & \text{subject to} \quad Ax = b \\ & \quad x_i \in \{0, 1\}, \quad i = 1 \dots n \end{aligned} \tag{726}$$

where  $\|x\|_0$  denotes cardinality of vector  $x$  (**a.k.a** 0-norm; not a convex function).

A minimal cardinality solution answers the question: “Which fewest linear combination of columns in  $A$  constructs vector  $b$ ?”. *Cardinality problems* have extraordinarily wide appeal, arising in many fields of science and across many disciplines. [355] [241] [195] [194] Yet designing an efficient algorithm to optimize cardinality has proved difficult. In this example, we also constrain the variable to be Boolean. The Boolean constraint forces an identical solution were the norm in problem (726) instead the 1-norm or 2-norm; *id est*, the two problems

$$\begin{aligned} & \underset{x}{\text{minimize}} \quad \|x\|_0 \\ (726) \quad & \text{subject to} \quad Ax = b \\ & \quad x_i \in \{0, 1\}, \quad i = 1 \dots n \end{aligned} \quad = \quad \begin{aligned} & \underset{x}{\text{minimize}} \quad \|x\|_1 \\ & \text{subject to} \quad Ax = b \\ & \quad x_i \in \{0, 1\}, \quad i = 1 \dots n \end{aligned} \tag{727}$$

are the same. The Boolean constraint makes the 1-norm problem nonconvex.

Given data

$$A = \begin{bmatrix} -1 & 1 & 8 & 1 & 1 & 0 \\ -3 & 2 & 8 & \frac{1}{2} & \frac{1}{3} & \frac{1}{2} - \frac{1}{3} \\ -9 & 4 & 8 & \frac{1}{4} & \frac{1}{9} & \frac{1}{4} - \frac{1}{9} \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{4} \end{bmatrix} \tag{728}$$

the obvious and desired solution to the problem posed,

$$x^* = e_4 \in \mathbb{R}^6 \tag{729}$$

has norm  $\|x^*\|_2 = 1$  and minimal cardinality; the minimum number of nonzero entries in vector  $x$ . The MATLAB **backslash command** `x=A\b`, for example, finds

$$x_M = \begin{bmatrix} \frac{2}{128} \\ 0 \\ \frac{5}{128} \\ 0 \\ \frac{90}{128} \\ 0 \end{bmatrix} \tag{730}$$

having norm  $\|x_M\|_2 = 0.7044$ . Coincidentally,  $x_M$  is a 1-norm solution; *id est*, an optimal solution to

$$\begin{aligned} & \underset{x}{\text{minimize}} \quad \|x\|_1 \\ & \text{subject to} \quad Ax = b \end{aligned} \tag{523}$$

The pseudoinverse solution (rounded)

$$x_p = A^\dagger b = \begin{bmatrix} -0.0456 \\ -0.1881 \\ 0.0623 \\ 0.2668 \\ 0.3770 \\ -0.1102 \end{bmatrix} \quad (731)$$

has least norm  $\|x_p\|_2 = 0.5165$ ; *id est*, the optimal solution to (§E.0.1.0.1)

$$\begin{array}{ll} \text{minimize}_x & \|x\|_2 \\ \text{subject to} & Ax = b \end{array} \quad (732)$$

Certainly none of the traditional methods provide  $x^* = e_4$  (729) because, and in general, for  $Ax = b$

$$\|\arg \inf \|x\|_2\|_2 \leq \|\arg \inf \|x\|_1\|_2 \leq \|\arg \inf \|x\|_0\|_2 \quad (733)$$

We can reformulate this minimal cardinality Boolean problem (726) as a semidefinite program: First transform the variable

$$x \triangleq (\hat{x} + \mathbf{1})_{\frac{1}{2}} \quad (734)$$

so  $\hat{x}_i \in \{-1, 1\}$ ; equivalently,

$$\begin{array}{ll} \text{minimize}_{\hat{x}} & \|(\hat{x} + \mathbf{1})_{\frac{1}{2}}\|_0 \\ \text{subject to} & A(\hat{x} + \mathbf{1})_{\frac{1}{2}} = b \\ & \delta(\hat{x}\hat{x}^T) = \mathbf{1} \end{array} \quad (735)$$

where  $\delta$  is the main-diagonal linear operator (§A.1). By assigning (§B.1)

$$G = \begin{bmatrix} \hat{x} \\ 1 \end{bmatrix} \begin{bmatrix} \hat{x}^T & 1 \end{bmatrix} = \begin{bmatrix} X & \hat{x} \\ \hat{x}^T & 1 \end{bmatrix} \triangleq \begin{bmatrix} \hat{x}\hat{x}^T & \hat{x} \\ \hat{x}^T & 1 \end{bmatrix} \in \mathbb{S}^{n+1} \quad (736)$$

problem (735) becomes equivalent to: (Theorem A.3.1.0.7)

$$\begin{array}{ll} \text{minimize}_{X \in \mathbb{S}^n, \hat{x} \in \mathbb{R}^n} & \mathbf{1}^T \hat{x} \\ \text{subject to} & A(\hat{x} + \mathbf{1})_{\frac{1}{2}} = b \\ & G = \begin{bmatrix} X & \hat{x} \\ \hat{x}^T & 1 \end{bmatrix} (\succeq 0) \\ & \delta(X) = \mathbf{1} \\ & \text{rank } G = 1 \end{array} \quad (737)$$

where solution is confined to rank-1 vertices of the *elliptope* in  $\mathbb{S}^{n+1}$  (§5.9.1.0.1) by the rank constraint, the positive semidefiniteness, and the equality constraints  $\delta(X) = \mathbf{1}$ . The rank constraint makes this problem nonconvex; by removing it<sup>4.15</sup> we get the semidefinite program

$$\begin{array}{ll} \text{minimize}_{X \in \mathbb{S}^n, \hat{x} \in \mathbb{R}^n} & \mathbf{1}^T \hat{x} \\ \text{subject to} & A(\hat{x} + \mathbf{1})_{\frac{1}{2}} = b \\ & G = \begin{bmatrix} X & \hat{x} \\ \hat{x}^T & 1 \end{bmatrix} \succeq 0 \\ & \delta(X) = \mathbf{1} \end{array} \quad (738)$$

---

<sup>4.15</sup>Relaxed problem (738) can also be derived via Lagrange duality; it is a dual of a dual program [*sic*] to (737). [340] [65, §5, exer.5.39] [439, §IV] [171, §11.3.4] The relaxed problem must therefore be convex having a larger feasible set; its optimal objective value represents a generally *loose* lower bound (1849) on the optimal objective of problem (737).

whose optimal solution  $x^*$  (734) is identical to that of minimal cardinality Boolean problem (726) if and only if  $\text{rank } G^* = 1$ .

Hope<sup>4.16</sup> of acquiring a rank-1 solution is not ill-founded because  $2^n$  ellipope vertices have rank 1 and because we are minimizing an affine function on a subset of the ellipope (Figure 155) containing rank-1 vertices; *id est*, by assumption that the feasible set of minimal cardinality Boolean problem (726) is nonempty, a desired solution resides on the ellipope relative boundary at a rank-1 vertex.<sup>4.17</sup>

For that data given in (728), our semidefinite program solver `sdpssol` [446] [447] (accurate in solution to approximately 1E-8)<sup>4.18</sup> finds optimal solution to (738)

$$\text{round}(G^*) = \begin{bmatrix} 1 & 1 & 1 & -1 & 1 & 1 & -1 \\ 1 & 1 & 1 & -1 & 1 & 1 & -1 \\ 1 & 1 & 1 & -1 & 1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 & -1 & 1 \\ 1 & 1 & 1 & -1 & 1 & 1 & -1 \\ 1 & 1 & 1 & -1 & 1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 & -1 & 1 \end{bmatrix} \quad (739)$$

near a rank-1 vertex of the ellipope in  $\mathbb{S}^{n+1}$  (Theorem 5.9.1.0.2); its sorted eigenvalues,

$$\lambda(G^*) = \begin{bmatrix} 6.99999977799099 \\ 0.00000022687241 \\ 0.00000002250296 \\ 0.00000000262974 \\ -0.00000000999738 \\ -0.00000000999875 \\ -0.00000001000000 \end{bmatrix} \quad (740)$$

Negative eigenvalues are undoubtedly finite-precision effects. Because the largest eigenvalue predominates by many orders of magnitude, we can expect to find a good approximation to a minimal cardinality Boolean solution by truncating all smaller eigenvalues. We find, indeed, the desired result (729)

$$x^* = \text{round} \left( \begin{bmatrix} 0.0000000127947 \\ 0.00000000527369 \\ 0.00000000181001 \\ 0.99999997469044 \\ 0.00000001408950 \\ 0.00000000482903 \end{bmatrix} \right) = e_4 \quad (741)$$

These numerical results are solver dependent; insofar, not all SDP solvers will return a rank-1 vertex solution.  $\square$

<sup>4.16</sup>A more deterministic approach to constraining rank and cardinality is in §4.7.0.0.12.

<sup>4.17</sup>Confinement to the ellipope can be regarded as a kind of normalization akin to matrix  $A$  column normalization suggested in [138] and explored in Example 4.2.3.1.2.

<sup>4.18</sup>A typically ignored limitation of interior-point solution methods is their relative accuracy of only about 1E-8 on a machine using 64-bit (*double precision*) floating-point arithmetic; *id est*, optimal solution  $x^*$  cannot be more accurate than square root of machine epsilon ( $\epsilon=2.2204E-16$ ). Nonzero primal-dual objective difference is not a good measure of solution accuracy.

**4.2.3.1.2 Example.** *Optimization over ellipope versus 1-norm polyhedron for minimal cardinality Boolean Example 4.2.3.1.1.*

A minimal cardinality problem is typically formulated via, what is by now, a standard practice [138] [75, §3.2, §3.4] of column normalization applied to a 1-norm problem surrogate like (523). Suppose we define a diagonal matrix

$$\Lambda \triangleq \begin{bmatrix} \|A(:, 1)\|_2 & & & & & \mathbf{0} \\ & \|A(:, 2)\|_2 & & & & \\ & & \ddots & & & \\ \mathbf{0} & & & & & \|A(:, 6)\|_2 \end{bmatrix} \in \mathbb{S}^6 \quad (742)$$

used to normalize the columns (assumed nonzero) of given noiseless data matrix  $A$ . Then approximate the minimal cardinality Boolean problem

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|x\|_0 \\ \text{subject to} & Ax = b \\ & x_i \in \{0, 1\}, \quad i=1 \dots n \end{array} \quad (726)$$

as

$$\begin{array}{ll} \underset{\tilde{y}}{\text{minimize}} & \|\tilde{y}\|_1 \\ \text{subject to} & A\Lambda^{-1}\tilde{y} = b \\ & \mathbf{1} \succeq \Lambda^{-1}\tilde{y} \succeq 0 \end{array} \quad (743)$$

where optimal solution

$$y^* = \text{round}(\Lambda^{-1}\tilde{y}^*) \quad (744)$$

The inequality in (743) relaxes Boolean constraint  $y_i \in \{0, 1\}$  from (726); bounding any solution  $y^*$  to a nonnegative unit hypercube whose vertices are binary numbers. Convex problem (743) is justified by the *convex envelope*

$$\text{cenv } \|x\|_0 \text{ on } \{x \in \mathbb{R}^n \mid \|x\|_\infty \leq \kappa\} = \frac{1}{\kappa} \|x\|_1 \quad (1520)$$

Donoho concurs with this particular formulation, equivalently expressible as a linear program via (519).

Approximation (743) is therefore equivalent to minimization of an affine function (§3.2) on a bounded polyhedron, whereas semidefinite program

$$\begin{array}{ll} \underset{X \in \mathbb{S}^n, \hat{x} \in \mathbb{R}^n}{\text{minimize}} & \mathbf{1}^T \hat{x} \\ \text{subject to} & A(\hat{x} + \mathbf{1})^{\frac{1}{2}} = b \\ & G = \begin{bmatrix} X & \hat{x} \\ \hat{x}^T & 1 \end{bmatrix} \succeq 0 \\ & \delta(X) = \mathbf{1} \end{array} \quad (738)$$

minimizes an affine function on an intersection of the ellipope with hyperplanes. Although the same Boolean solution is obtained from this approximation (743) as compared with semidefinite program (738), when given that particular data from Example 4.2.3.1.1, Singer confides a counterexample: Instead, given data

$$A = \begin{bmatrix} 1 & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & \frac{1}{\sqrt{2}} \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad (745)$$

then solving approximation (743) yields

$$y^* = \text{round} \left( \begin{bmatrix} 1 - \frac{1}{\sqrt{2}} \\ 1 - \frac{1}{\sqrt{2}} \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad (746)$$

(infeasible, with or without rounding, with respect to original problem (726)) whereas solving semidefinite program (738) produces

$$\text{round}(G^*) = \begin{bmatrix} 1 & 1 & -1 & 1 \\ 1 & 1 & -1 & 1 \\ -1 & -1 & 1 & -1 \\ 1 & 1 & -1 & 1 \end{bmatrix} \quad (747)$$

with sorted eigenvalues

$$\lambda(G^*) = \begin{bmatrix} 3.99999965057264 \\ 0.00000035942736 \\ -0.0000000000000000 \\ -0.00000001000000 \end{bmatrix} \quad (748)$$

Truncating all but the largest eigenvalue, from (734) we obtain (*confer*  $y^*$ )

$$x^* = \text{round}\left(\begin{bmatrix} 0.9999999625299 \\ 0.9999999625299 \\ 0.00000001434518 \end{bmatrix}\right) = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad (749)$$

the desired minimal cardinality Boolean result.  $\square$

#### 4.2.3.1.3 Exercise. Minimal cardinality Boolean art.

Assess general performance of standard-practice approximation (743) as compared with the proposed semidefinite program (738).  $\blacktriangledown$

#### 4.2.3.1.4 Exercise. Conic independence.

Matrix  $A$  from (728) is full-rank having three-dimensional nullspace. Find its four conically independent columns. (§2.10)<sup>4.19</sup> To what part of proper cone  $\mathcal{K} = \{Ax \mid x \succeq 0\}$  does vector  $b$  belong?  $\blacktriangledown$

#### 4.2.3.1.5 Exercise. Linear independence.

Show why wide matrix  $A$ , from compressed sensing problem (523) or (528), may be regarded full-rank without loss of generality. In other words: Is a minimal cardinality solution invariant to linear dependence of rows?  $\blacktriangledown$

## 4.3 Rank reduction

*... it is not clear generally how to predict rank  $X^*$  or rank  $S^*$  before solving the SDP problem.*

—Farid Alizadeh, 1995 [12, p.22]

The premise of rank reduction in semidefinite programming is: an optimal solution  $X^*$  found does not satisfy Barvinok's upper bound (275) on rank. The particular numerical algorithm solving a semidefinite program may have instead returned a high-rank optimal solution (§4.1.2; e.g., (707)) when a lower-rank optimal solution was expected. Rank reduction is a means to adjust rank of an optimal solution to (697P), returned by a solver, until it satisfies Barvinok's upper bound with the optimal objective value unchanged.

---

<sup>4.19</sup> Hint: §4.4.2.0.2, §4.6.2.0.2.

### 4.3.1 posit a perturbation of $X^*$

Recall (§4.1.2.1) that there is an extreme point of  $\mathcal{A} \cap \mathbb{S}_+^n$  satisfying upper bound (275) on rank. [25, §2.2] It is therefore sufficient to locate an extreme point of  $\mathcal{A} \cap \mathbb{S}_+^n$  whose primal objective value (697P) is optimal:<sup>4.20</sup> [126, §31.5.3] [268, §2.4] [269] [8, §3] [325]

Consider again affine subset

$$\mathcal{A} = \{X \in \mathbb{S}^n \mid A \text{ svec } X = b\} \quad (2273)$$

where for  $A_i \in \mathbb{S}^n$

$$A = \begin{bmatrix} \text{svec}(A_1)^T \\ \vdots \\ \text{svec}(A_m)^T \end{bmatrix} \in \mathbb{R}^{m \times n(n+1)/2} \quad (698)$$

Given any optimal solution  $X^*$  to SDP

$$\begin{array}{ll} \underset{X \in \mathbb{S}^n}{\text{minimize}} & \langle C, X \rangle \\ \text{subject to} & X \in \mathcal{A} \cap \mathbb{S}_+^n \end{array} \quad (697P)$$

whose rank does not satisfy upper bound (275), we posit existence of a set of perturbations

$$\{t_j B_j \mid t_j \in \mathbb{R}, B_j \in \mathbb{S}^n, j=1 \dots n\} \quad (750)$$

to  $X^*$  such that, for some  $0 \leq i \leq n$  and scalars  $\{t_j, j=1 \dots i\}$ ,

$$X^* + \sum_{j=1}^i t_j B_j \quad (751)$$

becomes an extreme point of  $\mathcal{A} \cap \mathbb{S}_+^n$  and remains an optimal solution to (697P). Membership of (751) to affine subset  $\mathcal{A}$  is secured, for the  $i^{\text{th}}$  perturbation, by demanding

$$\langle B_i, A_j \rangle = 0, \quad j=1 \dots m \quad (752)$$

while membership to positive semidefinite cone  $\mathbb{S}_+^n$  is insured by small perturbation (761). Feasibility of (751) is certified in this manner, whereas optimality is proved in §4.3.3.

The following simple algorithm has low computational intensity and locates an optimal extreme point, assuming nontrivial solution: given optimal primal solution  $X^*$

#### 4.3.1.0.1 Procedure. Rank reduction.

[426]

```

initialize:  $B_i = \mathbf{0} \quad \forall i$ 
for iteration  $i=1 \dots n$ 
{
  1. compute a nonzero perturbation matrix  $B_i$  (755) of  $X^* + \sum_{j=1}^{i-1} t_j^* B_j$ 
  2. maximize  $t_i$  (761)
  subject to  $X^* + \sum_{j=1}^{i-1} t_j^* B_j + t_i B_i \in \mathbb{S}_+^n$ 
}

```

¶

A rank-reduced optimal solution is then

$$X^* \leftarrow X^* + \sum_{j=1}^i t_j^* B_j \quad (753)$$

---

<sup>4.20</sup> There is no known construction for Barvinok's tighter result (280).

—Monique Laurent, 2004

### 4.3.2 rank perturbation form

Perturbations of  $X^*$  are independent of constants  $C \in \mathbb{S}^n$  and  $b \in \mathbb{R}^m$  in primal and dual problems (697). Numerical accuracy of any rank-reduced result, found by perturbation of an initial optimal solution  $X^*$ , is therefore quite dependent upon initial accuracy of  $X^*$ .

#### 4.3.2.0.1 Definition. Matrix step function.

(confer §A.6.3.2.1)

Define the signum-like quasiconcave real function  $\psi : \mathbb{S}^n \rightarrow \mathbb{R}$

$$\psi(Z) \triangleq \begin{cases} 1, & Z \succeq 0 \\ -1, & \text{otherwise} \end{cases} \quad (754)$$

The value  $-1$  is taken for indefinite or nonzero negative semidefinite argument.<sup>4.21</sup>  $\triangle$

Deza & Laurent [126, §31.5.3] prove: every perturbation matrix  $B_i$ ,  $i=1 \dots n$ , is of the form

$$B_i = -\psi(Z_i)R_iZ_iR_i^T \in \mathbb{S}^n \quad (755)$$

where

$$X^* \triangleq R_1R_1^T, \quad X^* + \sum_{j=1}^{i-1} t_j^* B_j \triangleq X_i = R_iR_i^T \in \mathbb{S}^n \quad (756)$$

where the optimal  $t_j^*$  are scalars and  $R_i \in \mathbb{R}^{n \times \rho}$  is full-rank and thin where

$$\rho \triangleq \text{rank} \left( X^* + \sum_{j=1}^{i-1} t_j^* B_j \right) = \text{rank } X_i \quad (757)$$

and where  $Z_i \in \mathbb{S}^\rho$  is found at each iteration  $i$  by solving a simple feasibility problem:<sup>4.22</sup>

$$\begin{array}{ll} \text{find} & R_iZ_iR_i^T \neq \mathbf{0} \\ \text{subject to} & \langle Z_i, R_i^T A_j R_i \rangle = 0, \quad j=1 \dots m \end{array} \quad (758)$$

Were there a sparsity pattern common to each member of set  $\{R_i^T A_j R_i \in \mathbb{S}^\rho, j=1 \dots m\}$ , then a good choice for  $Z_i$  has 1 in each entry corresponding to a 0 in the pattern; *id est*, a sparsity pattern complement. At iteration  $i$

$$X^* + \sum_{j=1}^{i-1} t_j^* B_j + t_i B_i = R_i(I - t_i \psi(Z_i)Z_i)R_i^T \quad (759)$$

By fact (1604), therefore

$$X^* + \sum_{j=1}^{i-1} t_j^* B_j + t_i B_i \succeq 0 \Leftrightarrow \mathbf{1} - t_i \psi(Z_i)\lambda(Z_i) \succeq 0 \quad (760)$$

where  $\lambda(Z_i) \in \mathbb{R}^\rho$  denotes the eigenvalues of  $Z_i$ . Necessity and sufficiency are due to the facts:  $R_i$  can be completed to a nonsingular matrix (§A.3.1.0.5.c), and  $I - t_i \psi(Z_i)Z_i$  can

<sup>4.21</sup>Because of how  $\mathbf{0}$  and indefinites are handled,  $\psi$  is not an odd function; *id est*,  $\psi(-Z) \neq -\psi(Z)$ .

<sup>4.22</sup>A simple method of solution is closed-form projection of a nonzero random point  $Z_i$  on that proper subspace of isometrically isomorphic  $\mathbb{R}^{\rho(\rho+1)/2}$  specified by the constraints. (§E.5.0.0.7) Such a solution is nontrivial assuming the specified intersection of hyperplanes is not the origin; guaranteed by  $\rho(\rho+1)/2 > m$ . This geometric intuition, about forming a perturbation, is indeed what bounds any solution's rank from below;  $m$  is fixed by the number of equality constraints in (697P) while rank  $\rho$  decreases with each iteration  $i$ . Otherwise, we might iterate indefinitely.

be padded with zeros while maintaining equivalence in (759). Maximization of each  $t_i$ , in step 2 of Procedure 4.3.1.0.1, reduces rank of (759) so locates a new point on the boundary  $\partial(\mathcal{A} \cap \mathbb{S}_+^n)$ .<sup>4.23</sup> Maximization of  $t_i$  thereby has closed form;

$$(t_i^*)^{-1} = \max \{\psi(Z_i)\lambda(Z_i)_k, k=1 \dots \rho\} \quad (761)$$

When  $Z_i$  is indefinite, direction of perturbation (determined by  $\psi(Z_i)$ ) is arbitrary. We may take an early exit, from the Procedure, were all feasible  $R_i Z_i R_i^T$  to become  $\{\mathbf{0}\}$  or were  $\rho$  to become equal to 1 (assuming a nontrivial solution) or were

$$\text{rank} [\text{svec}(R_i^T A_1 R_i) \text{ svec}(R_i^T A_2 R_i) \cdots \text{svec}(R_i^T A_m R_i)] = \rho(\rho + 1)/2 \quad (762)$$

(277) which characterizes rank  $\rho$  of any [sic] extreme point in  $\mathcal{A} \cap \mathbb{S}_+^n$ . [268, §2.4] [269]

**Proof.** Assuming the form of every perturbation matrix is indeed (755), then by (758)

$$\text{svec } Z_i \perp [\text{svec}(R_i^T A_1 R_i) \text{ svec}(R_i^T A_2 R_i) \cdots \text{svec}(R_i^T A_m R_i)] \quad (763)$$

By orthogonal complement we have

$$\begin{aligned} & \text{rank} [\text{svec}(R_i^T A_1 R_i) \text{ svec}(R_i^T A_2 R_i) \cdots \text{svec}(R_i^T A_m R_i)]^\perp \\ & + \text{rank} [\text{svec}(R_i^T A_1 R_i) \text{ svec}(R_i^T A_2 R_i) \cdots \text{svec}(R_i^T A_m R_i)] = \rho(\rho + 1)/2 \end{aligned} \quad (764)$$

When  $Z_i$  can only be  $\mathbf{0}$ , then the perturbation is null because an extreme point has been found; thus

$$[\text{svec}(R_i^T A_1 R_i) \text{ svec}(R_i^T A_2 R_i) \cdots \text{svec}(R_i^T A_m R_i)]^\perp = \mathbf{0} \quad (765)$$

from which the stated result (762) directly follows. ♦

### 4.3.3 Optimality of perturbed $X^*$

We show that the optimal objective value is unaltered by perturbation (755); *id est*,

$$\langle C, X^* + \sum_{j=1}^i t_j^* B_j \rangle = \langle C, X^* \rangle \quad (766)$$

**Proof.** From Corollary 4.2.3.0.1 we have the necessary and sufficient relationship between optimal primal and dual solutions under assumption of nonempty primal feasible cone interior  $\mathcal{A} \cap \text{intr } \mathbb{S}_+^n$ :

$$S^* X^* = S^* R_1 R_1^T = X^* S^* = R_1 R_1^T S^* = \mathbf{0} \quad (767)$$

This means  $\mathcal{R}(R_1) \subseteq \mathcal{N}(S^*)$  and  $\mathcal{R}(S^*) \subseteq \mathcal{N}(R_1^T)$ . From (756) and (759), after 0-padding  $Z_i$  for dimensional compatibility, come the sequence:

---

<sup>4.23</sup>This holds because rank of a positive semidefinite matrix in  $\mathbb{S}^n$  is diminished below  $n$  by the number of its 0 eigenvalues (1614), and because a positive semidefinite matrix having one or more 0 eigenvalues corresponds to a point on the PSD cone boundary (196).

$$\begin{aligned}
X^* &= R_1 R_1^T \\
X^* + t_1^* B_1 &= R_2 R_2^T = R_1(I - t_1^* \psi(Z_1) Z_1) R_1^T \\
X^* + t_1^* B_1 + t_2^* B_2 &= R_3 R_3^T = R_2(I - t_2^* \psi(Z_2) Z_2) R_2^T = R_1 \sqrt{I - t_1^* \psi(Z_1) Z_1} (I - t_2^* \psi(Z_2) Z_2) \sqrt{I - t_1^* \psi(Z_1) Z_1} R_1^T \\
&\vdots \\
X^* + \sum_{j=1}^i t_j^* B_j &= R_1 \left( \prod_{j=1}^i \sqrt{I - t_j^* \psi(Z_j) Z_j} \right) \left( \prod_{j=i}^1 \sqrt{I - t_j^* \psi(Z_j) Z_j} \right) R_1^T, \tag{768} \quad i > 0
\end{aligned}$$

where second product counts backwards. Substituting  $C = \text{svec}^{-1}(A^T y^*) + S^*$  from (697),

$$\begin{aligned}
\langle C, X^* + \sum_{j=1}^i t_j^* B_j \rangle &= \left\langle \text{svec}^{-1}(A^T y^*) + S^*, R_1 \prod_{j=1}^i \sqrt{I - t_j^* \psi(Z_j) Z_j} \prod_{j=i}^1 \sqrt{I - t_j^* \psi(Z_j) Z_j} R_1^T \right\rangle \\
&= \left\langle \sum_{k=1}^m y_k^* A_k, X^* + \sum_{j=1}^i t_j^* B_j \right\rangle \\
&= \left\langle \sum_{k=1}^m y_k^* A_k + S^*, X^* \right\rangle = \langle C, X^* \rangle \tag{769}
\end{aligned}$$

because  $\langle B_i, A_j \rangle = 0 \quad \forall i, j$  by design (752).  $\spadesuit$

#### 4.3.3.0.1 Example. $A \delta(X) = b$ .

This academic example demonstrates that a solution found by rank reduction can certainly have rank less than Barvinok's upper bound (275): Assume that a given vector  $b$  belongs to the conic hull of columns of a given matrix  $A$

$$A = \begin{bmatrix} -1 & 1 & 8 & 1 & 1 \\ -3 & 2 & 8 & \frac{1}{2} & \frac{1}{3} \\ -9 & 4 & 8 & \frac{1}{4} & \frac{1}{9} \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad b = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{4} \end{bmatrix} \in \mathbb{R}^m \tag{770}$$

Consider the convex optimization problem

$$\begin{aligned}
&\underset{X \in \mathbb{S}^5}{\text{minimize}} \quad \text{tr } X \\
&\text{subject to} \quad X \succeq 0 \\
&\quad A \delta(X) = b
\end{aligned} \tag{771}$$

that minimizes the 1-norm of the main diagonal; *id est*, problem (771) is the same as

$$\begin{aligned}
&\underset{X \in \mathbb{S}^5}{\text{minimize}} \quad \|\delta(X)\|_1 \\
&\text{subject to} \quad X \succeq 0 \\
&\quad A \delta(X) = b
\end{aligned} \tag{772}$$

that finds a solution to  $A \delta(X) = b$ . Rank-3 solution  $X^* = \delta(x_M)$  is optimal, where (*confer*(730))

$$x_M = \begin{bmatrix} \frac{2}{128} \\ 0 \\ \frac{5}{128} \\ 0 \\ \frac{90}{128} \end{bmatrix} \tag{773}$$

Yet upper bound (275) predicts existence of at most a

$$\text{rank-}\left(\left\lfloor \frac{\sqrt{8m+1}-1}{2} \right\rfloor = 2\right) \quad (774)$$

feasible solution from  $m=3$  equality constraints. To find a lower rank  $\rho$  optimal solution to (771) (barring combinatorics), we invoke Procedure 4.3.1.0.1:

**Initialize:**  $C=I$ ,  $\rho=3$ ,  $A_j \triangleq \delta(A(j,:))$ ,  $j=1, 2, 3$ ,  $X^*=\delta(x_M)$ ,  $m=3$ ,  $n=5$ .  
 $\{$

Iteration  $i=1$ :

$$\text{Step 1: } R_1 = \begin{bmatrix} \sqrt{\frac{2}{128}} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & \sqrt{\frac{5}{128}} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \sqrt{\frac{90}{128}} \end{bmatrix}.$$

$$\begin{array}{ll} \text{find}_{Z_1 \in \mathbb{S}^3} & R_1 Z_1 R_1^T \neq \mathbf{0} \\ \text{subject to} & \langle Z_1, R_1^T A_j R_1 \rangle = 0, \quad j=1, 2, 3 \end{array} \quad (775)$$

A nonzero randomly selected matrix  $Z_1$ , having  $\mathbf{0}$  main diagonal, is a solution yielding nonzero perturbation matrix  $B_1$ . Choose arbitrarily

$$Z_1 = \mathbf{1}\mathbf{1}^T - I \in \mathbb{S}^3 \quad (776)$$

Then (rounding)

$$B_1 = \begin{bmatrix} 0 & 0 & 0.0247 & 0 & 0.1048 \\ 0 & 0 & 0 & 0 & 0 \\ 0.0247 & 0 & 0 & 0 & 0.1657 \\ 0 & 0 & 0 & 0 & 0 \\ 0.1048 & 0 & 0.1657 & 0 & 0 \end{bmatrix} \quad (777)$$

**Step 2:**  $t_1^* = 1$  because  $\lambda(Z_1) = [-1 \ -1 \ 2]^T$ . So,

$$X^* \leftarrow \delta(x_M) + t_1^* B_1 = \begin{bmatrix} \frac{2}{128} & 0 & 0.0247 & 0 & 0.1048 \\ 0 & 0 & 0 & 0 & 0 \\ 0.0247 & 0 & \frac{5}{128} & 0 & 0.1657 \\ 0 & 0 & 0 & 0 & 0 \\ 0.1048 & 0 & 0.1657 & 0 & \frac{90}{128} \end{bmatrix} \quad (778)$$

has rank  $\rho \leftarrow 1$  and produces the same optimal objective value.

}

□

#### 4.3.3.0.2 Exercise. Rank reduction of maximal complementarity.

Apply rank reduction Procedure 4.3.1.0.1 to the *maximal complementarity example* (§4.1.2.3.1). Demonstrate a rank-1 solution; which can certainly be found (by Barvinok's Proposition 2.9.3.0.1) because there is only one equality constraint. ▼

### 4.3.4 thoughts regarding rank reduction

Because rank reduction Procedure 4.3.1.0.1 is guaranteed only to produce another optimal solution conforming to Barvinok's upper bound (275), the Procedure will not necessarily produce solutions of arbitrarily low rank; but if they exist, the Procedure can. Arbitrariness of search direction, when matrix  $Z_i$  becomes indefinite (mentioned on page 241), and the enormity of choices for  $Z_i$  (758) are liabilities for this algorithm.

#### 4.3.4.1 inequality constraints

The question naturally arises: what to do when a semidefinite program (not in prototypical form (697))<sup>4.24</sup> has linear inequality constraints of the form

$$\alpha_i^T \text{svec } X \preceq \varphi_i, \quad i = 1 \dots k \quad (779)$$

where  $\{\varphi_i\}$  are given scalars and  $\{\alpha_i\}$  are given vectors. One expedient way to handle this circumstance is to convert the inequality constraints to equality constraints by introducing a slack variable  $\gamma$ ; *id est*,

$$\alpha_i^T \text{svec } X + \gamma_i = \varphi_i, \quad i = 1 \dots k, \quad \gamma \succeq 0 \quad (780)$$

thereby converting the problem to prototypical form.

Alternatively, we say the  $i^{\text{th}}$  inequality constraint is *active* when it is met with equality; *id est*, when for particular  $i$  in (779),  $\alpha_i^T \text{svec } X^* = \varphi_i$ . An optimal high-rank solution  $X^*$  is, of course, feasible (satisfying all the constraints). But for the purpose of rank reduction, inactive inequality constraints are ignored while active inequality constraints are interpreted as equality constraints. In other words, we take the union of active inequality constraints (as equalities) with equality constraints  $A \text{svec } X = b$  to form a composite affine subset  $\hat{\mathcal{A}}$  substituting for (2273). Then we proceed with rank reduction of  $X^*$  as though the semidefinite program were in prototypical form (697P).

## 4.4 Cardinality reduction

Analogous to rank reduction of semidefinite variable in SDP (§4.3), cardinality reduction of vector variable in LP means: to lower cardinality of an optimal solution to (696p) (found by numerical solver) while leaving the optimal objective value unchanged.

### 4.4.1 perturbation of $x^*$

Given affine subset

$$\mathcal{A} = \{x \in \mathbb{R}^n \mid Ax = b\} \quad (151)$$

where

$$A = \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} \in \mathbb{R}^{m \times n} \quad (150)$$

and given any optimal solution  $x^*$  to LP

$$\begin{aligned} & \underset{x}{\text{minimize}} && c^T x \\ & \text{subject to} && x \succeq 0 \\ & && Ax = b \end{aligned} \quad (696\text{p})$$

---

<sup>4.24</sup>Contemporary numerical packages for solving semidefinite programs can solve a range of problems wider than prototype (697). Generally, they do so by transforming a given problem into prototypical form by introducing new constraints and variables. [12] [447] We are momentarily considering a departure from the primal prototype that augments the constraint set with linear inequalities.

whose cardinality is not minimal, an extreme point of  $\mathcal{A} \cap \mathbb{R}_+^n$  (whose primal objective value (696p) is optimal) would possess reduced cardinality. To reveal such an extreme point, we posit existence of a set of perturbations to  $x^*$  (like those in §4.3.1)

$$\{t_j \beta_j \mid t_j \in \mathbb{R}, \beta_j \in \mathbb{R}^n, j=1 \dots n\} \quad (781)$$

such that, for some  $0 \leq i \leq n$  and set of scalars  $\{t_j, j=1 \dots i\}$ ,

$$x^* + \sum_{j=1}^i t_j \beta_j \quad (782)$$

becomes extreme and optimal. Membership of (782) to affine subset  $\mathcal{A}$  is guaranteed, for the  $i^{\text{th}}$  perturbation, by constraints

$$\langle \beta_i, a_j \rangle = 0, \quad j=1 \dots m \quad (783)$$

while membership to nonnegative orthant  $\mathbb{R}_+^n$  is insured by small perturbation (792). Thus, feasibility of (782) is certain.

#### 4.4.2 cardinality perturbation form

Perturbation of  $x^*$  is independent of vector constants  $c \in \mathbb{R}^n$  and  $b \in \mathbb{R}^m$  in primal and dual problems (696). Every perturbation  $\beta_i, i=1 \dots n$ , is a vector of the form

$$\beta_i = -\psi(\delta(z_i)) z_i \circ x_i \in \mathbb{R}^n \quad (784)$$

where

$$x^* \triangleq x_1, \quad x^* + \sum_{j=1}^{i-1} t_j^* \beta_j \triangleq x_i \in \mathbb{R}^n \quad (785)$$

where the optimal  $t_j^*$  are scalars and where  $z_i$  is found at each iteration  $i$  by solving a simple feasibility problem:

$$\begin{aligned} & \underset{z_i \in \mathbb{R}^n}{\text{find}} && z_i \circ x_i \neq \mathbf{0} \\ & \text{subject to} && \langle z_i, a_j \circ x_i \rangle = 0, \quad j=1 \dots m \end{aligned} \quad (786)$$

Cardinality  $\rho$  of  $x_i \in \mathbb{R}^n$  is equivalent to number of its nonzero entries:

$$\rho \triangleq \text{card} \left( x^* + \sum_{j=1}^{i-1} t_j^* \beta_j \right) = \text{card } x_i \quad (787)$$

At iteration  $i$

$$x^* + \sum_{j=1}^{i-1} t_j^* \beta_j + t_i \beta_i = (\mathbf{1} - t_i \psi(\delta(z_i)) z_i) \circ x_i \quad (788)$$

Hence, the sequence

$$\begin{aligned} x^* &= x_1 \\ x^* + t_1^* \beta_1 &= x_2 = (\mathbf{1} - t_1 \psi(\delta(z_1)) z_1) \circ x_1 \\ x^* + t_1^* \beta_1 + t_2^* \beta_2 &= x_3 = (\mathbf{1} - t_2 \psi(\delta(z_2)) z_2) \circ x_2 = (\mathbf{1} - t_1 \psi(\delta(z_1)) z_1) \circ (\mathbf{1} - t_2 \psi(\delta(z_2)) z_2) \circ x_1 \\ &\vdots \\ x^* + \sum_{j=1}^i t_j^* \beta_j &= \left( \prod_{j=1}^i \delta(\mathbf{1} - t_j \psi(\delta(z_j)) z_j) \right) x_1, \end{aligned} \quad (789) \quad i > 0$$

from which it follows (in order of iteration):

$$x^* + \sum_{j=1}^{i-1} t_j^* \beta_j + t_i \beta_i \succeq 0 \Leftrightarrow \mathbf{1} - t_i \psi(\delta(z_i)) z_i \succeq 0, \quad i = 1 \dots n \quad (790)$$

The following algorithm locates an optimal extreme point, assuming nontrivial solution: given any optimal primal solution  $x^*$

#### 4.4.2.0.1 Procedure. Cardinality reduction.

```

initialize:  $\beta_i = \mathbf{0}$   $\forall i$ 
for iteration  $i = 1 \dots n$ 
{
    1. compute a nonzero perturbation vector  $\beta_i$  (784) of  $x^* + \sum_{j=1}^{i-1} t_j^* \beta_j$ 
    2. maximize  $t_i$  (792)
        subject to  $x^* + \sum_{j=1}^{i-1} t_j^* \beta_j + t_i \beta_i \succeq 0$ 
}

```

¶

A cardinality-reduced optimal solution is then

$$x^* \leftarrow x^* + \sum_{j=1}^i t_j^* \beta_j \quad (791)$$

Maximization of  $t_i$ , in step 2 of Procedure 4.4.2.0.1, reduces cardinality of (788) so locates a new point on boundary  $\partial(\mathcal{A} \cap \mathbb{R}_+^n)$ . Maximization of  $t_i$  thereby has closed form;

$$(t_i^*)^{-1} = \max \{\psi(\delta(z_i))z_i(k), k = 1 \dots n\} \quad (792)$$

We may exit early, from the Procedure, were all feasible  $z_i \circ x_i$  to become  $\{\mathbf{0}\}$  or were cardinality  $\rho$  to become 1 or were

$$\text{rank}[a_1 \circ x_i \ a_2 \circ x_i \ \dots \ a_m \circ x_i] = \rho \quad (793)$$

which characterizes cardinality  $\rho$  of any extreme point in  $\mathcal{A} \cap \mathbb{R}_+^n$ .

#### 4.4.2.0.2 Example. $Ax = b$ .

Cardinality minimization is often at odds with norm minimization because these two objectives can compete; e.g., §4.2.3.1.1. Yet, prior knowledge of optimal norm objective value may facilitate a cardinality minimization problem. If optimal solution  $x^*$  were known to be binary with particular cardinality  $\rho$ , for example, then a linear constraint on the variable  $\mathbf{1}^T x = \rho$  might be warranted because  $\rho = \|x\|_1$  for a binary variable. Columns of this particular  $A$  matrix

$$A = \begin{bmatrix} -1 & 1 & 8 & 1 & 1 & 0 \\ -3 & 2 & 8 & \frac{1}{2} & \frac{1}{3} & \frac{1}{2} - \frac{1}{3} \\ -9 & 4 & 8 & \frac{1}{4} & \frac{1}{9} & \frac{1}{4} - \frac{1}{9} \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad b = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{4} \end{bmatrix} \in \mathbb{R}^m \quad (728)$$

constitute generators of a pointed polyhedral cone  $\mathcal{K}$ .<sup>4.25</sup> Vector  $b$  predetermines optimal 1-norm, of a binary variable, to be 1 or 2. This convex feasibility problem

$$\begin{aligned}
 & \text{find } x \in \mathbb{R}^6 \\
 & \text{subject to } x \succeq 0 \\
 & \quad Ax = b \\
 & \quad c^T x = 1
 \end{aligned} \quad (794)$$

---

<sup>4.25</sup> Columns {1,5,2,6} are c.i. generators, {1,5} {5,2} {2,6} {6,1} generate facets, {3,4} are interior to  $\mathcal{K}$ .

brings objective  $c^T x$  ( $c = \mathbf{1}$ ) down into the constraints. Were cardinality-1 solution found, feasible  $x$  would certainly be binary. Because minimization of  $c^T x$  is forgone, conditions for 0-duality gap (308) are unmet; objective value cannot be maintained as in §4.3.3.

$$x_G = \begin{bmatrix} \frac{2}{159} \\ 0 \\ \frac{5}{159} \\ 0 \\ \frac{121}{159} \\ \frac{31}{159} \end{bmatrix} \quad (795)$$

Cardinality-4  $x_G$  solves (794). Ignoring norm constraint  $c^T x = 1$ , Procedure 4.4.2.0.1 may be invoked to find a lesser cardinality solution:

Initialize:  $c = \mathbf{1}$ ,  $\rho = 1$ ,  $a_j$ ,  $j = 1, 2, 3$  (150)(p.244),  $x^* = x_G$ ,  $m = 3$ ,  $n = 6$ .  
 {

Iteration  $i = 1$ :

Step 1:  $x_1 = x^*$ .

$$\begin{array}{ll} \text{find} & z_1 \circ x_1 \neq \mathbf{0} \\ z_1 \in \mathbb{R}^6 & \text{subject to } \langle z_1, a_j \circ x_1 \rangle = 0, \quad j = 1, 2, 3 \end{array} \quad (796)$$

Choose

$$z_1 = \left[ -\frac{159}{128} \quad 1 \quad -\frac{159}{128} \quad 1 \quad \frac{1546}{3963} \quad \frac{159}{31} \right]^T \quad (797)$$

Then (784)

$$\beta_1 = \left[ -\frac{1}{64} \quad 0 \quad -\frac{5}{128} \quad 0 \quad \frac{19}{64} \quad 1 \right]^T \quad (798)$$

Step 2:  $t_1^* = \frac{128}{159}$ . So,

$$x^* \leftarrow x_G + t_1^* \beta_1 = [0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 1]^T \quad (799)$$

has cardinality  $\rho \leftarrow 2$ .

}

Further iterations  $i$  produce  $z_i = \mathbf{0}$ . □

As illustrated by Example 4.4.2.0.2, cardinality reduction can fail (at (784)) to find a minimal cardinality solution when  $x_1$  has a 0-entry in a minimal cardinality location. This result instigates search for a new method:

## 4.5 Rank constraint by Convex Iteration

We generalize the trace heuristic (§7.2.2.1), for finding low-rank optimal solutions to semidefinite programs of a more general form:

### 4.5.1 constraining rank of semidefinite matrices

Consider a *semidefinite feasibility problem* of the form

$$\begin{array}{ll} \text{find} & G \\ G \in \mathbb{S}^N & \text{subject to } G \in \mathcal{C} \\ & G \succeq 0 \\ & \text{rank } G \leq n \end{array} \quad (800)$$

where  $\mathcal{C}$  is a convex set presumed to contain positive semidefinite matrices of rank  $n$  or less; *id est*,  $\mathcal{C}$  intersects the positive semidefinite cone boundary. We propose: this

rank-constrained feasibility problem can be equivalently expressed as iteration of the convex problem sequence (801) and (1872a):

$$\begin{aligned} & \underset{G \in \mathbb{S}^N}{\text{minimize}} \quad \langle G, W \rangle \\ & \text{subject to} \quad G \in \mathcal{C} \\ & \quad G \succeq 0 \end{aligned} \tag{801}$$

where *direction vector*<sup>4.26</sup>  $W \in \mathbb{S}^N$  is an optimal solution to the following semidefinite program, for  $0 \leq n \leq N - 1$

$$\begin{aligned} \sum_{i=n+1}^N \lambda(G^*)_i &= \underset{W \in \mathbb{S}^N}{\text{minimize}} \quad \langle G^*, W \rangle \\ &\text{subject to} \quad 0 \preceq W \preceq I \\ &\quad \text{tr } W = N - n \end{aligned} \tag{1872a}$$

whose feasible set is a Fantope ([§2.3.2.0.1](#)),<sup>4.27</sup> and where  $G^*$  is an optimal solution to problem (801) given some iterate  $W$ . The idea is to iterate solution of (801) and (1872a) until convergence as defined in [§4.5.1.2](#): (*confer* (837))

$$\sum_{i=n+1}^N \lambda(G^*)_i = \langle G^*, W^* \rangle = \lambda(G^*)^T \lambda(W^*) \triangleq 0 \tag{802}$$

defines *global optimality* of the iteration; a vanishing objective that is a certificate of global optimality but cannot be guaranteed. Inner product of eigenvalues follows from (1748) and properties of commutative matrix products ([p.494](#)). *Optimal direction vector*  $W^*$  is defined as any positive semidefinite matrix yielding optimal solution  $G^*$  of rank  $n$  or less to then convex equivalent (801) of feasibility problem (800):

$$\begin{aligned} & \underset{G \in \mathbb{S}^N}{\text{find}} \quad G \\ (800) \quad & \text{subject to} \quad G \in \mathcal{C} \quad \equiv \quad \underset{G \in \mathbb{S}^N}{\text{minimize}} \quad \langle G, W^* \rangle \\ & \quad G \succeq 0 \\ & \quad \text{rank } G \leq n \quad \text{subject to} \quad G \in \mathcal{C} \\ & \quad G \succeq 0 \end{aligned} \tag{801}$$

*id est*, any direction vector for which the last  $N - n$  nonincreasingly ordered eigenvalues  $\lambda$  of  $G^*$  are zero.

In any semidefinite feasibility problem, a solution of least rank must be an extreme point of the feasible set.<sup>4.28</sup> This means there exists a hyperplane supporting the feasible set at that extreme point. ([§2.11](#)) Then there must exist a linear objective function such that this least-rank feasible solution optimizes the resultant semidefinite program.

We emphasize that convex problem (801) is not a relaxation of rank-constrained feasibility problem (800); at global optimality, convex iteration (801) (1872a) makes it instead an *equivalent problem*.

#### 4.5.1.1 direction matrix interpretation

(*confer* [§4.6.1.2](#)) The feasible set of direction matrices in (1872a) is the convex hull of outer product of all rank- $(N - n)$  orthonormal matrices; *videlicet*,

$$\text{conv} \left\{ UU^T \mid U \in \mathbb{R}^{N \times N-n}, U^T U = I \right\} = \left\{ A \in \mathbb{S}^N \mid I \succeq A \succeq 0, \langle I, A \rangle = N - n \right\} \tag{91}$$

<sup>4.26</sup>Search *direction*  $W$  is a hyperplane-normal pointing opposite to direction of movement describing minimization of a real linear function  $\langle G, W \rangle$  ([p.62](#)).

<sup>4.27</sup>Sum of eigenvalues follows from a result of Ky Fan ([p.533](#)).

<sup>4.28</sup>which follows by *extremes theorem* [2.8.1.1.1](#), by rank of a sum of positive semidefinite matrices ([1620](#)) ([262](#)), and by definition of extreme point ([170](#)) for which no convex combination can produce it: If a least rank solution were expressible as a convex combination of feasible points, then there could exist feasible matrices of lesser rank.

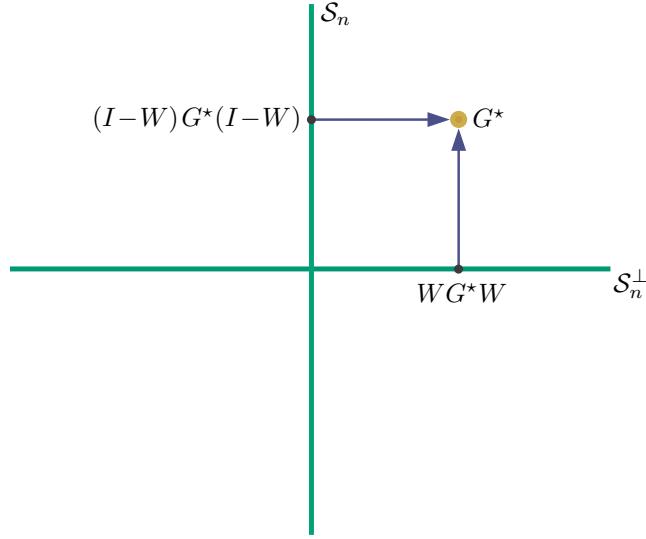


Figure 95: (*confer* Figure 192) Projection of  $G^*$  on subspace  $\mathcal{S}_n$  of rank  $\leq n$  matrices whose nullspace contains  $\mathcal{N}(G^*)$ . This direction  $W$  is closed-form solution to (1872a).

This set (93), argument to  $\text{conv}\{\}$ , comprises the extreme points of this Fantope (91). An optimal solution  $W$  to (1872a), that is an extreme point, is known in closed form (p.533): Given ordered diagonalization  $G^* = Q\Lambda Q^T \in \mathbb{S}_+^N$  (§A.5.1), then direction matrix  $W = U^*U^{*T}$  is optimal and extreme where  $U^* = Q(:, n+1:N) \in \mathbb{R}^{N \times N-n}$ . Eigenvalue vector  $\lambda(W)$  has 1 in each entry corresponding to the  $N-n$  smallest entries of  $\delta(\Lambda)$  and has 0 elsewhere. By (225) (228), polar direction  $-W$  can be regarded as pointing toward the set of all rank- $n$  (or less) positive semidefinite matrices whose nullspace contains that of  $G^*$ . For that particular closed-form solution  $W$ , consequent to Theobald (p.493), (*confer*(839))

$$\sum_{i=n+1}^N \lambda(G^*)_i = \langle G^*, W \rangle = \lambda(G^*)^T \lambda(W) \geq 0 \quad (803)$$

This is the connection to cardinality minimization of vectors;<sup>4.29</sup> *id est*, eigenvalue  $\lambda$  cardinality (rank) is analogous to vector  $x$  cardinality via (839): for positive semidefinite  $X$

$$\begin{aligned} \sum_i \lambda(X)_i &= \text{tr } X &= \|X\|_2^* &\Leftrightarrow \|x\|_1 \\ \sqrt{\sum_i \lambda(X)_i^2} &= \sqrt{\text{tr } X^2} &= \|X\|_{\text{F}} &\Leftrightarrow \|x\|_2 \\ \max_i \{\lambda(X)_i\} & &= \|X\|_2 &\Leftrightarrow \|x\|_{\infty} \end{aligned} \quad (804)$$

So that this method, for constraining rank, will not be misconstrued under closed-form solution  $W$  to (1872a): Define (*confer*(225))

$$\mathcal{S}_n \triangleq \{(I-W)G(I-W) \mid G \in \mathbb{S}^N\} = \{X \in \mathbb{S}^N \mid \mathcal{N}(X) \supseteq \mathcal{N}(G^*)\} \quad (805)$$

<sup>4.29</sup> not trace minimization of a nonnegative diagonal matrix  $\delta(x)$  as in [153, §1] [336, §2]. To make rank-constrained problem (800) resemble cardinality problem (535), we could make  $\mathcal{C}$  an affine subset:

$$\begin{array}{ll} \text{find} & X \in \mathbb{S}^N \\ \text{subject to} & A \text{svec } X = b \\ & X \succeq 0 \\ & \text{rank } X \leq n \end{array}$$

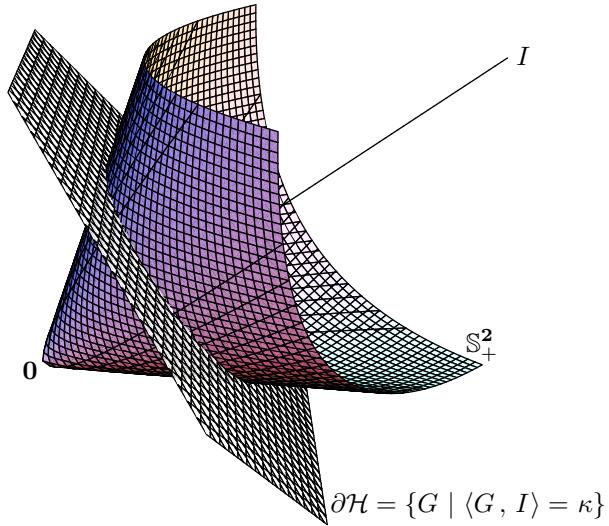


Figure 96: (confer Figure 112) Trace heuristic can be interpreted as minimization of a hyperplane, with normal  $I$ , over positive semidefinite cone drawn here in isometrically isomorphic  $\mathbb{R}^3$ . Polar of direction vector  $W = I$  points toward origin.

as the symmetric subspace of rank  $\leq n$  matrices whose nullspace contains  $\mathcal{N}(G^*)$ . Then projection of  $G^*$  on  $\mathcal{S}_n$  is  $(I - W)G^*(I - W)$ . (§E.7) Direction of projection is  $-WG^*W$ . (Figure 95)  $\text{tr}(WG^*W)$  is a measure of proximity to  $\mathcal{S}_n$  because its orthogonal complement is  $\mathcal{S}_n^\perp = \{WGW \mid G \in \mathbb{S}^N\}$ ; the point being, convex iteration (incorporating constrained  $\text{tr}(WGW) = \langle G, W \rangle$  minimization) is not a projection method: certainly, not on these two subspaces. Proposed convex iteration is neither *dual projection* (Figure 191) or *alternating projection* (Figure 195).

Closed-form solution  $W$  to problem (1872a), though efficient, comes with a *caveat*: there exist cases where this projection matrix solution  $W$  does not provide the shortest route to an optimal rank- $n$  solution  $G^*$ ; *id est*, direction  $W$  is not unique. So we sometimes choose to solve (1872a) instead of employing a known closed-form solution.

When direction matrix  $W = I$ , as in the trace heuristic for example, then  $-W$  points directly at the origin (the rank-0 PSD matrix, Figure 96). Vector inner-product of an optimization variable with direction matrix  $W$  is therefore a generalization of the trace heuristic (§7.2.2.1) for rank minimization;  $-W$  is instead trained toward the boundary of the positive semidefinite cone.

#### 4.5.1.2 convergence

We study convergence to ascertain conditions under which a direction matrix will reveal a feasible solution  $G$ , of rank  $n$  or less, to semidefinite program (801). Denote by  $W^*$  a particular optimal direction matrix from semidefinite program (1872a) such that (802) holds (feasible rank  $G \leq n$  found). Then we define *global optimality* of the iteration (801) (1872a) to correspond with this vanishing vector inner-product (802) of optimal solutions.

Because this iterative technique for constraining rank is not a projection method, it can find a rank- $n$  solution  $G^*$  ((802) will be satisfied) only if at least one exists in the feasible set of program (801).

**4.5.1.2.1 Proof.** Suppose  $\langle G^*, W \rangle = \tau$  is satisfied for some nonnegative constant  $\tau$  after any particular iteration (801) (1872a) of the two minimization problems. Once a particular value of  $\tau$  is achieved, it can never be exceeded by subsequent iterations because existence of feasible  $G$  and  $W$  having that vector inner-product  $\tau$  has been established simultaneously in each problem. Because the infimum of vector inner-product of two positive semidefinite matrix variables is zero, the nonincreasing sequence of iterations is thus bounded below hence convergent because any bounded monotonic sequence in  $\mathbb{R}$  is convergent. [289, §1.2] [43, §1.1] *Local optimality* to some nonnegative objective value  $\tau$  is thereby established. ♦

*Local optimality*, in this context, means convergence of  $\langle G^*, W \rangle$  to a *fixed point* of possibly infeasible rank. Only local optimality can be established because objective  $\langle G, W \rangle$ , when instead regarded simultaneously in two variables  $(G, W)$ , is generally multimodal. (§3.14.0.0.3)

Local optimality, convergence to  $\langle G^*, W \rangle = \tau \neq 0$  and definition of a *stall*, never implies nonexistence of a rank- $n$  feasible solution to (801). Conversely, a nonexistent rank- $n$  feasible solution would mean certain failure ( $\tau \neq 0$ ) to achieve global optimality by definition (802). But, as proved, convex iteration always converges to a local optimum if not a global one.

When a rank- $n$  feasible solution to (801) exists, it remains an open problem to state conditions under which  $\langle G^*, W^* \rangle = \tau = 0$  (802) is achieved by iterative solution of semidefinite programs (801) and (1872a). Then rank  $G^* \leq n$  and pair  $(G^*, W^*)$  becomes a globally optimal fixed point of iteration. There can be no proof of convergence to a global optimum because of the implicit high-dimensional multimodal manifold in variables  $(G, W)$ . When stall occurs, direction vector  $W$  can be manipulated to steer out; *e.g.*, reversal of search direction as in Example 4.7.0.0.1, or reinitialization to a random rank- $(N-n)$  matrix in the same PSD cone face (§2.9.2.3) demanded by the current iterate: given ordered diagonalization  $G^* = Q\Lambda Q^T \in \mathbb{S}^N$ , then  $W = U^* \Phi U^{*T}$  where  $U^* = Q(:, n+1:N) \in \mathbb{R}^{N \times N-n}$  and where eigenvalue vector  $\lambda(W)_{1:N-n} = \lambda(\Phi)$  has nonnegative uniformly distributed random entries in  $(0, 1]$  by selection of  $\Phi \in \mathbb{S}_+^{N-n}$  while  $\lambda(W)_{N-n+1:N} = \mathbf{0}$ . Zero eigenvalues act as memory while randomness largely reduces likelihood of stall. When this direction works, rank and objective sequence  $\langle G^*, W \rangle$  (with respect to iteration) tend to be noisily monotonic.

**4.5.1.2.2 Exercise.** Completely positive semidefinite matrix. [41]

Given rank-2 positive semidefinite matrix  $G = \begin{bmatrix} 0.50 & 0.55 & 0.20 \\ 0.55 & 0.61 & 0.22 \\ 0.20 & 0.22 & 0.08 \end{bmatrix}$ , find a positive factorization  $G = X^T X$  (1042) by solving

$$\begin{aligned} & \underset{X \in \mathbb{R}^{2 \times 3}}{\text{find}} \quad X \geq \mathbf{0} \\ & \text{subject to} \quad Z = \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} \succeq 0 \\ & \text{rank } Z \leq 2 \end{aligned} \tag{806}$$

via convex iteration. ▼

**4.5.1.2.3 Exercise.** Nonnegative matrix factorization.

Given rank-2 nonnegative matrix  $X = \begin{bmatrix} 17 & 28 & 42 \\ 16 & 47 & 51 \\ 17 & 82 & 72 \end{bmatrix}$ , find a nonnegative factorization

$$X = WH \tag{807}$$

by solving

$$\begin{array}{ll}
 \underset{A \in \mathbb{S}^3, B \in \mathbb{S}^3, W \in \mathbb{R}^{3 \times 2}, H \in \mathbb{R}^{2 \times 3}}{\text{find}} & W, H \\
 \text{subject to} & Z = \begin{bmatrix} I & W^T & H \\ W & A & X \\ H^T & X^T & B \end{bmatrix} \succeq 0 \\
 & W \geq \mathbf{0} \\
 & H \geq \mathbf{0} \\
 & \text{rank } Z \leq 2
 \end{array} \tag{808}$$

which follows from the fact, at optimality,

$$Z^* = \begin{bmatrix} I \\ W \\ H^T \end{bmatrix} [I \ W^T \ H] \tag{809}$$

Use the known closed-form solution for a direction vector  $Y$  to regulate rank by convex iteration; set  $Z^* = Q\Lambda Q^T \in \mathbb{S}^8$  to an ordered diagonalization and  $U^* = Q(:, 3:8) \in \mathbb{R}^{8 \times 6}$ , then  $Y = U^*U^{*\top}$  (§4.5.1.1).

In summary, initialize  $Y$  then iterate numerical solution of (convex) semidefinite program

$$\begin{array}{ll}
 \underset{A \in \mathbb{S}^3, B \in \mathbb{S}^3, W \in \mathbb{R}^{3 \times 2}, H \in \mathbb{R}^{2 \times 3}}{\text{minimize}} & \langle Z, Y \rangle \\
 \text{subject to} & Z = \begin{bmatrix} I & W^T & H \\ W & A & X \\ H^T & X^T & B \end{bmatrix} \succeq 0 \\
 & W \geq \mathbf{0} \\
 & H \geq \mathbf{0}
 \end{array} \tag{810}$$

with  $Y = U^*U^{*\top}$  until convergence (which is to a global optimum, and occurs in very few iterations for this instance).  $\blacktriangledown$

Now, an application to optimal regulation of affine dimension:

#### 4.5.1.2.4 Example. Sensor-Network Localization and Wireless Location.

Heuristic solution to a sensor-network localization problem, proposed by Carter, Jin, Saunders, & Ye in [79],<sup>4.30</sup> is limited to two Euclidean dimensions and applies semidefinite programming (SDP) to little subproblems. There, a large network is partitioned into smaller subnetworks (as small as one *sensor* - a mobile point, whereabouts unknown) and then semidefinite programming and heuristics called SPASELOC are applied to localize each and every partition by two-dimensional distance geometry. Their partitioning procedure is one-pass, yet termed *iterative*; a term applicable only insofar as adjoining partitions can share localized sensors and *anchors* (absolute sensor positions known *a priori*). But there is no iteration on the entire network, hence the term “iterative” is perhaps inappropriate. As partitions are selected based on “rule sets” (heuristics, not geographics), they also term the partitioning *adaptive*. But no adaptation of a partition actually occurs once it has been determined.

---

<sup>4.30</sup>The paper constitutes [Jin's dissertation](#) for University of Toronto [243] although her name appears as second author. Ye's authorship is honorary.

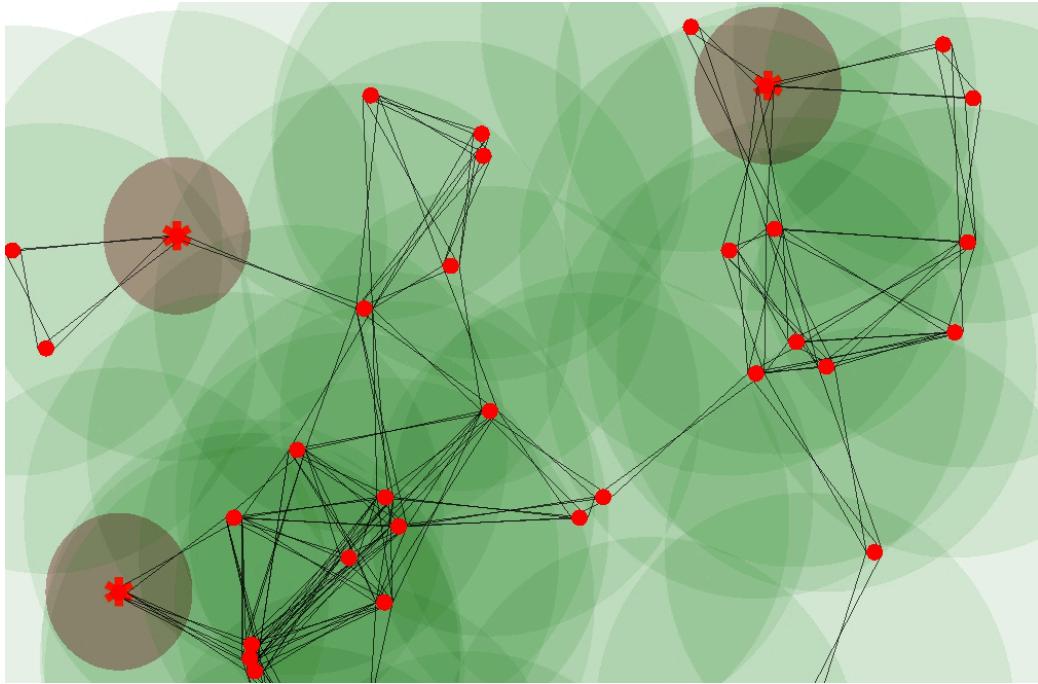


Figure 97: Sensor-network localization in  $\mathbb{R}^2$ , illustrating connectivity and circular radio-range per *sensor*. Smaller dark grey regions each hold an *anchor* at their center; known fixed sensor positions. Sensor/anchor distance is measurable with negligible uncertainty for sensor within those grey regions. (Graphic by [Geoff Merrett](#).)

One can reasonably argue that semidefinite programming methods are unnecessary for localization of small partitions of large sensor networks. [310] [94] In the past, these nonlinear localization problems were solved algebraically and computed by least squares solution to hyperbolic equations; called *multilateration*.<sup>4.31</sup> [256] [297] Indeed, practical contemporary numerical methods for global positioning (GPS) by satellite do not rely on convex optimization. [323]

Modern distance geometry is inextricably melded with semidefinite programming. The beauty of semidefinite programming, as relates to localization, lies in convex expression of classical multilateration: So & Ye showed [357] that the problem of finding unique solution, to a noiseless nonlinear system describing the common point of intersection of hyperspheres in real Euclidean vector space, can be expressed as a semidefinite program via distance geometry.

But the need for SDP methods in Carter & Jin *et alii* is enigmatic for two more reasons: 1) guessing solution to a partition whose intersensor measurement data or connectivity is inadequate for localization by distance geometry, 2) reliance on complicated and extensive heuristics for partitioning a large network that could instead be efficiently solved whole by one semidefinite program [251, §3]. While partitions range in size between 2 and 10 sensors, 5 sensors optimally, heuristics provided are only for two spatial dimensions (no higher-dimensional heuristics are proposed). For these small numbers it remains unclarified as to precisely what advantage is gained over traditional

---

<sup>4.31</sup>**Multilateration** - literally, *having many sides*; shape of a geometric figure formed by nearly intersecting lines of position. In navigation systems, therefore: Obtaining a *fix* from multiple lines of position. Multilateration can be regarded as noisy trilateration.

least squares: it is difficult to determine what part of their noise performance is attributable to SDP and what part is attributable to their heuristic geometry.

Partitioning of large sensor networks is a compromise to rapid growth of SDP computational intensity with problem size. But when impact of noise on distance measurement is of most concern, one is averse to a partitioning scheme because noise-effects vary inversely with problem size. [56, §2.2] (§5.13.2) Since an individual partition's solution is not iterated in Carter & Jin and is interdependent with adjoining partitions, we expect errors to propagate from one partition to the next; the ultimate partition solved, expected to suffer most.

Heuristics often fail on real-world data because of unanticipated circumstances. When heuristics fail, generally they are repaired by adding more heuristics. Tenuous is any presumption, for example, that distance measurement errors have distribution characterized by circular contours of equal probability about an unknown sensor-location. (Figure 97) That presumption effectively appears within Carter & Jin's optimization problem statement as affine equality constraints relating unknowns to distance measurements that are corrupted by noise. Yet in most all urban environments, this measurement noise is more aptly characterized by ellipsoids of varying orientation and eccentricity as one recedes from a sensor. (Figure 151) Each unknown sensor must therefore instead be bound to its own particular range of distance, primarily determined by the terrain.<sup>4.32</sup> The nonconvex problem we must instead solve is:

$$\begin{aligned} \underset{i,j \in \mathcal{I}}{\text{find}} \quad & \{x_i, x_j\} \\ \text{subject to} \quad & \underline{d}_{ij} \leq \|x_i - x_j\|^2 \leq \overline{d}_{ij} \end{aligned} \tag{811}$$

where  $x_i$  represents sensor location, and where  $\underline{d}_{ij}$  and  $\overline{d}_{ij}$  respectively represent lower and upper bounds on measured distance-square from  $i^{\text{th}}$  to  $j^{\text{th}}$  sensor (or from sensor to anchor). Figure 102 illustrates contours of equal sensor-location uncertainty. By establishing these individual upper and lower bounds, orientation and eccentricity can effectively be incorporated into the problem statement.

Generally speaking, there can be no unique solution to the sensor-network localization problem because there is no unique formulation; that is the art of Optimization. Any optimal solution obtained depends on whether or how a network is partitioned, whether distance data is complete, presence of noise, and how the problem is formulated. When a particular formulation is a convex optimization problem, then the set of all optimal solutions forms a convex set containing the actual or true localization. Measurement noise precludes equality constraints representing distance. The optimal solution set is consequently expanded; necessitated by introduction of distance inequalities admitting more and higher-rank solutions. Even were the optimal solution set a single point, it is not necessarily the true localization because there is little hope of exact localization by any algorithm once significant noise is introduced.

Carter & Jin gauge performance of their heuristics to the SDP formulation of author Biswas whom they regard as vanguard to the art. [15, §1] Biswas posed localization as an optimization problem minimizing a distance measure. [50] [48] Intuitively, minimization of any distance measure yields compacted solutions; (*confer* §6.7.0.0.1) precisely the anomaly motivating Carter & Jin. Their two-dimensional heuristics outperformed Biswas' localizations both in execution-time and proximity to the desired result. Perhaps, instead of heuristics, Biswas' approach to localization can be improved: [47] [49].

---

<sup>4.32</sup>A distinct contour map corresponding to each anchor is required in practice.

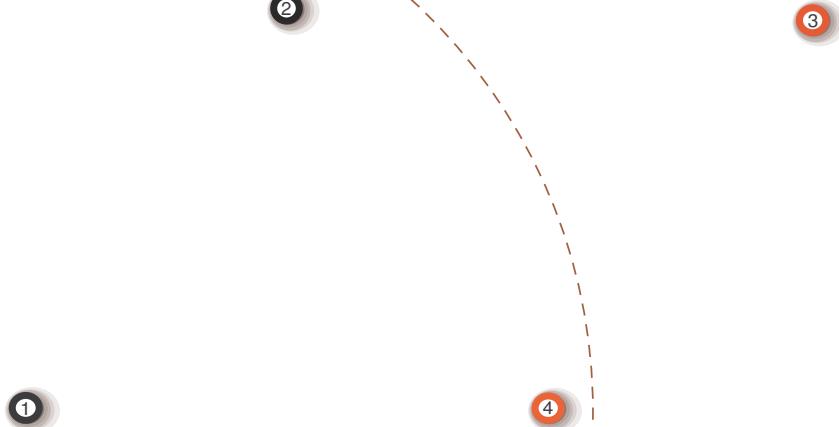


Figure 98: 2-lattice in  $\mathbb{R}^2$ , hand-drawn. Nodes 3 and 4 are anchors; remaining nodes are sensors. Radio range of sensor 1 indicated by arc.

The sensor-network localization problem is considered difficult. [15, §2] Rank constraints in optimization are considered more difficult. Control of affine dimension in Carter & Jin is suboptimal because of implicit projection on  $\mathbb{R}^2$ . In what follows, we present the localization problem as a semidefinite program (equivalent to (811)) having an explicit rank constraint which controls affine dimension of an optimal solution. We show how to achieve that rank constraint only if the feasible set contains a matrix of desired rank. Our problem formulation is extensible to any spatial dimension.

### proposed standardized test

Jin proposes an academic test in two-dimensional real Euclidean space  $\mathbb{R}^2$  that we adopt. In essence, this test is a localization of sensors and anchors arranged in a regular triangular lattice. Lattice connectivity is solely determined by sensor radio range; a connectivity graph is assumed incomplete. In the interest of test standardization, we propose adoption of a few small examples: Figure 98 through Figure 101 and their particular connectivity represented by matrices (812) through (815) respectively.

$$\begin{matrix} 0 & \bullet & ? & \bullet \\ \bullet & 0 & \bullet & \bullet \\ ? & \bullet & 0 & \circ \\ \bullet & \bullet & \circ & 0 \end{matrix} \quad (812)$$

Matrix entries *dot*  $\bullet$  indicate measurable distance between *nodes* while unknown distance is denoted by  $?$  (*question mark*). Matrix entries *hollow dot*  $\circ$  represent known distance between anchors (to high accuracy) while zero distance is denoted 0. Because measured distances are quite unreliable in practice, our solution to the localization problem substitutes a distinct range of possible distance for each measurable distance; equality constraints exist only for anchors.

Anchors are chosen so as to increase difficulty for algorithms dependent on existence of sensors in their convex hull. The challenge is to find a solution in two dimensions close to the true sensor positions given incomplete noisy intersensor distance information.

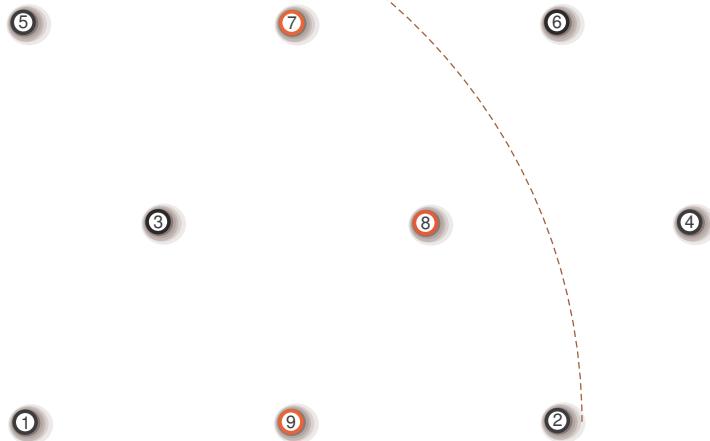


Figure 99: 3-lattice in  $\mathbb{R}^2$ , hand-drawn. Nodes 7, 8, and 9 are anchors; remaining nodes are sensors. Radio range of sensor 1 indicated by arc.

$$\begin{matrix}
 0 & \bullet & \bullet & ? & \bullet & ? & ? & \bullet & \bullet \\
 \bullet & 0 & \bullet & \bullet & ? & \bullet & ? & \bullet & \bullet \\
 \bullet & \bullet & 0 & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\
 ? & \bullet & \bullet & 0 & ? & \bullet & \bullet & \bullet & \bullet \\
 \bullet & ? & \bullet & ? & 0 & \bullet & \bullet & \bullet & \bullet \\
 ? & \bullet & \bullet & \bullet & \bullet & 0 & \bullet & \bullet & \bullet \\
 ? & ? & \bullet & \bullet & \bullet & \bullet & 0 & \circ & \circ \\
 \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \circ & 0 & \circ \\
 \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \circ & \circ & 0
 \end{matrix} \tag{813}$$

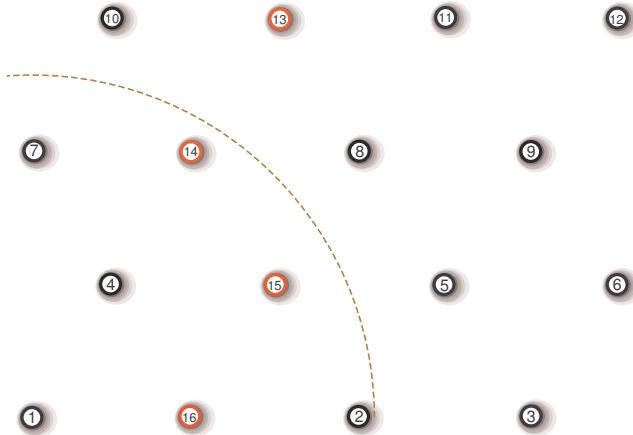


Figure 100: 4-lattice in  $\mathbb{R}^2$ , hand-drawn. Nodes 13, 14, 15, and 16 are anchors; remaining nodes are sensors. Radio range of sensor 1 indicated by arc.

0	?	?	●	?	?	●	?	?	?	?	?	?	?	?	●	●
?	0	●	●	●	●	?	●	?	?	?	?	?	?	?	●	●
?	●	0	?	●	●	?	●	?	?	?	?	?	?	?	●	●
●	●	?	0	●	?	●	●	?	●	?	?	?	●	●	●	●
?	●	●	●	●	0	●	?	●	●	?	●	●	●	●	●	●
?	●	●	?	●	0	?	●	●	?	●	●	?	?	?	?	?
●	?	?	●	?	?	0	?	?	●	?	?	●	●	●	●	●
?	●	?	●	●	●	?	0	●	●	●	●	●	●	●	●	●
?	?	●	?	●	●	?	●	0	?	●	●	?	●	?	●	?
?	?	?	●	?	?	●	●	?	0	●	?	●	●	●	?	?
?	?	?	?	●	●	?	●	●	●	0	●	●	●	●	●	?
?	?	?	?	●	●	?	●	●	?	●	0	?	?	?	?	?
?	?	?	●	●	?	●	●	●	●	●	?	0	○	○	○	○
?	●	?	●	●	?	●	●	?	●	●	?	○	0	○	○	○
●	●	●	●	●	?	●	●	●	●	●	?	○	○	○	○	○
●	●	●	●	●	?	●	●	?	?	?	?	○	○	○	○	0

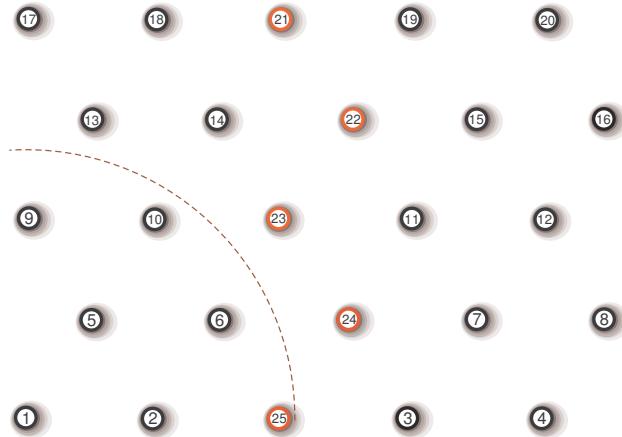


Figure 101: 5-lattice in  $\mathbb{R}^2$ . Nodes 21 through 25 are anchors.

•	?	?	•	•	?	?	•	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	
•	0	?	?	•	•	?	?	•	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	•	•
?	?	0	•	?	•	•	•	?	?	•	•	?	?	?	?	?	?	?	?	?	?	?	?	?	?
•	•	?	?	0	•	?	?	•	•	?	?	•	•	?	?	?	?	?	?	?	?	?	?	?	?
•	•	•	?	•	0	•	?	•	•	•	?	?	•	?	?	?	?	?	?	?	?	?	?	?	?
?	?	•	•	?	•	0	•	?	?	•	•	?	?	•	•	?	?	?	?	?	?	?	?	?	?
?	?	•	•	?	•	0	?	?	•	•	?	?	•	•	?	?	?	?	?	?	?	?	?	?	?
•	?	?	•	•	?	?	0	•	?	?	•	•	?	?	?	•	•	?	?	?	?	?	?	?	?
?	•	?	?	•	•	?	?	•	0	•	?	?	•	•	?	?	?	•	•	•	•	•	•	•	•
?	?	•	•	?	?	•	0	•	?	?	•	•	?	?	?	•	•	?	•	•	•	•	•	•	?
?	?	?	?	•	?	?	•	0	•	?	?	0	•	?	?	•	•	?	?	•	•	?	?	?	?
?	?	?	?	•	•	?	?	•	•	•	?	•	0	•	?	?	•	•	•	•	•	•	•	•	?
?	?	?	?	?	•	•	?	?	•	•	?	?	•	0	•	?	?	•	•	•	?	•	?	?	?
?	?	?	?	?	•	?	?	•	?	?	•	•	?	?	0	•	?	?	•	?	?	•	?	?	?
?	?	?	?	?	?	•	?	?	?	•	•	?	?	•	?	?	•	0	•	?	•	•	?	?	?
?	?	?	?	?	?	?	•	?	?	•	•	?	?	•	?	?	•	0	•	•	•	?	?	?	?
?	?	?	?	?	?	?	?	•	?	?	•	•	?	?	•	?	•	•	•	0	o	o	o	o	
?	?	?	?	?	?	?	?	?	•	?	?	•	•	•	•	•	•	0	o	o	o	o	o	o	
?	?	?	?	?	?	?	?	?	?	•	•	•	•	?	•	?	•	•	?	o	o	o	o	o	
?	?	?	?	?	?	?	?	?	?	?	•	•	•	?	•	?	•	•	?	o	o	o	o	o	
?	?	?	?	?	?	?	?	?	?	?	?	•	•	?	?	?	?	?	?	o	o	o	o	o	
?	?	?	?	?	?	?	?	?	?	?	?	?	•	•	?	?	?	?	?	o	o	o	o	o	
?	?	?	?	?	?	?	?	?	?	?	?	?	?	•	•	?	?	?	?	o	o	o	o	o	
?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	•	•	?	?	?	o	o	o	o	o	
?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	•	•	?	?	?	o	o	o	o	o

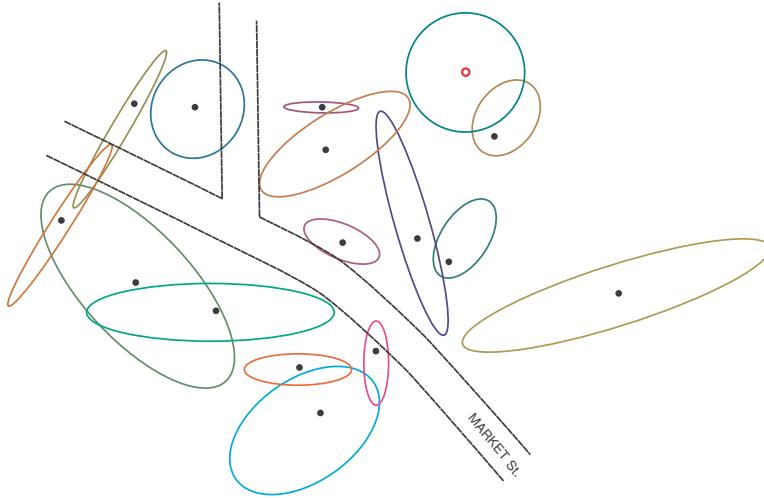


Figure 102: Location uncertainty ellipsoid in  $\mathbb{R}^2$  for each of 15 sensors  $\bullet$  within three city blocks in downtown San Francisco. (Data by [Polaris Wireless](#).)

### problem statement

Ascribe points in a list  $\{x_\ell \in \mathbb{R}^n, \ell = 1 \dots N\}$  to the columns of a matrix  $X$ ;

$$X = [x_1 \ \cdots \ x_N] \in \mathbb{R}^{n \times N} \quad (77)$$

where  $N$  is regarded as cardinality of list  $X$ . Positive semidefinite matrix  $X^T X$ , formed from inner product of the list, is a *Gram matrix*; [280, §3.6]

$$G = X^T X = \begin{bmatrix} \|x_1\|^2 & x_1^T x_2 & x_1^T x_3 & \cdots & x_1^T x_N \\ x_2^T x_1 & \|x_2\|^2 & x_2^T x_3 & \cdots & x_2^T x_N \\ x_3^T x_1 & x_3^T x_2 & \|x_3\|^2 & \ddots & x_3^T x_N \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ x_N^T x_1 & x_N^T x_2 & x_N^T x_3 & \cdots & \|x_N\|^2 \end{bmatrix} \in \mathbb{S}_+^N \quad (1042)$$

where  $\mathbb{S}_+^N$  is the convex cone of  $N \times N$  positive semidefinite matrices in the symmetric matrix subspace  $\mathbb{S}^N$ .

Existence of noise precludes measured distance from the input data. We instead assign measured distance to a range estimate specified by individual upper and lower bounds:  $\bar{d}_{ij}$  is an upper bound on distance-square from  $i^{\text{th}}$  to  $j^{\text{th}}$  sensor, while  $\underline{d}_{ij}$  is a lower bound. These bounds become the input data. Each measurement range is presumed different from the others because of measurement uncertainty; *e.g.*, Figure 102.

Our mathematical treatment of anchors and sensors is not dichotomized.<sup>4.33</sup> A sensor position that is known *a priori* to high accuracy (with absolute certainty)  $\hat{x}_i$  is called an *anchor*. Then the sensor-network localization problem (811) can be expressed equivalently: Given a number  $m$  of anchors and a set of indices  $\mathcal{I}$  (corresponding to all measurable distances  $\bullet$ ), for  $0 < n < N$

<sup>4.33</sup> Wireless location problem thus stated identically; difference being: fewer sensors.

$$\begin{aligned}
& \underset{G \in \mathbb{S}^N, X \in \mathbb{R}^{n \times N}}{\text{find}} && X \\
\text{subject to} & \quad d_{ij} \leq \langle G, (e_i - e_j)(e_i - e_j)^T \rangle \leq \bar{d}_{ij} \quad \forall (i, j) \in \mathcal{I} \\
& \langle G, e_i e_i^T \rangle &= \|\check{x}_i\|^2, \quad i = N-m+1 \dots N \\
& \langle G, (e_i e_j^T + e_j e_i^T)/2 \rangle &= \check{x}_i^T \check{x}_j, \quad i < j, \quad \forall i, j \in \{N-m+1 \dots N\} \\
& X(:, N-m+1:N) &= [\check{x}_{N-m+1} \ \cdots \ \check{x}_N] \\
& Z = \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} &\succeq 0 \\
& \text{rank } Z &= n
\end{aligned} \tag{816}$$

where  $e_i$  is the  $i^{\text{th}}$  member of the standard basis for  $\mathbb{R}^N$ . Distance-square

$$d_{ij} = \|x_i - x_j\|_2^2 = \langle x_i - x_j, x_i - x_j \rangle \tag{1029}$$

is related to Gram matrix entries  $G \triangleq [g_{ij}]$  by vector inner-product

$$\begin{aligned}
d_{ij} &= g_{ii} + g_{jj} - 2g_{ij} \\
&= \langle G, (e_i - e_j)(e_i - e_j)^T \rangle = \text{tr}(G^T(e_i - e_j)(e_i - e_j)^T)
\end{aligned} \tag{1044}$$

hence the scalar inequalities. Each linear equality constraint in  $G \in \mathbb{S}^N$  represents a hyperplane in isometrically isomorphic Euclidean vector space  $\mathbb{R}^{N(N+1)/2}$ , while each linear inequality pair represents a convex Euclidean body known as *slab*.<sup>4.34</sup> By Schur complement (§A.4), any solution  $(G, X)$  provides comparison with respect to the positive semidefinite cone

$$G \succeq X^T X \tag{1082}$$

which is a convex relaxation of the desired equality constraint

$$\begin{bmatrix} I & X \\ X^T & G \end{bmatrix} = \begin{bmatrix} I & X \\ X^T & \end{bmatrix} [I \ X] \tag{1083}$$

The rank constraint insures this equality holds, by Theorem A.4.0.1.3, thus restricting solution to  $\mathbb{R}^n$ . Assuming full-rank solution (list)  $X$

$$\text{rank } Z = \text{rank } G = \text{rank } X \tag{817}$$

### convex equivalent problem statement

Problem statement (816) is nonconvex because of the rank constraint. We do not eliminate or ignore the rank constraint; rather, we find a convex way to enforce it: for  $0 < n < N$

$$\begin{aligned}
& \underset{G \in \mathbb{S}^N, X \in \mathbb{R}^{n \times N}}{\text{minimize}} && \langle Z, W \rangle \\
\text{subject to} & \quad d_{ij} \leq \langle G, (e_i - e_j)(e_i - e_j)^T \rangle \leq \bar{d}_{ij} \quad \forall (i, j) \in \mathcal{I} \\
& \langle G, e_i e_i^T \rangle &= \|\check{x}_i\|^2, \quad i = N-m+1 \dots N \\
& \langle G, (e_i e_j^T + e_j e_i^T)/2 \rangle &= \check{x}_i^T \check{x}_j, \quad i < j, \quad \forall i, j \in \{N-m+1 \dots N\} \\
& X(:, N-m+1:N) &= [\check{x}_{N-m+1} \ \cdots \ \check{x}_N] \\
& Z = \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} &\succeq 0
\end{aligned} \tag{818}$$

---

<sup>4.34</sup> an intersection of two parallel but opposing halfspaces (Figure 13). In terms of position  $X$ , this distance slab can be thought of as a thick *hypershell* instead of a hypersphere boundary.

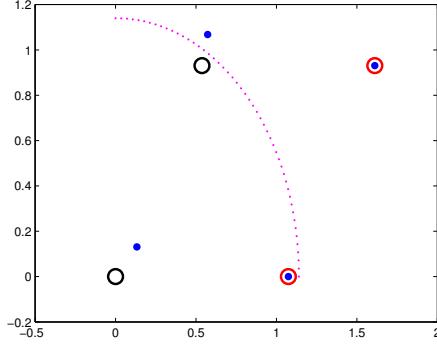


Figure 103: Typical solution for 2-lattice in Figure 98 with noise factor  $\eta = 0.1$  (821). Two red rightmost nodes are anchors; two remaining nodes are sensors. Radio range of sensor 1 indicated by arc; radius = 1.14. Actual sensor indicated by target  $\circ$  while its localization is indicated by bullet  $\bullet$ . Rank-2 solution found in 1 iteration (818) (1872a) subject to reflection error.

Convex function  $\text{tr } Z$  is a well-known heuristic whose sole purpose is to represent convex envelope of rank  $Z$ . (§7.2.2.1) In this convex optimization problem (818), a semidefinite program, we substitute a vector inner-product objective function for trace;

$$\text{tr } Z = \langle Z, I \rangle \leftarrow \langle Z, W \rangle \quad (819)$$

a generalization of the trace heuristic for minimizing convex envelope of rank, where  $W \in \mathbb{S}_+^{N+n}$  is constant with respect to (818). Matrix  $W$  is normal to a hyperplane in  $\mathbb{S}^{N+n}$  minimized over a convex feasible set specified by the constraints in (818). Matrix  $W$  is chosen so  $-W$  points in direction of rank- $n$  feasible solutions  $G$ . For properly chosen  $W$ , problem (818) becomes an equivalent to (816). Thus the purpose of vector inner-product objective (819) is to locate a rank- $n$  feasible Gram matrix assumed existent on the boundary of positive semidefinite cone  $\mathbb{S}_+^N$ , as explained beginning in §4.5.1; how to choose direction vector  $W$  is explained there and in what follows:

#### direction matrix $W$

Denote by  $Z^*$  an optimal composite matrix from semidefinite program (818). Then for  $Z^* \in \mathbb{S}^{N+n}$  whose eigenvalues  $\lambda(Z^*) \in \mathbb{R}^{N+n}$  are arranged in nonincreasing order, (Ky Fan)

$$\begin{aligned} \sum_{i=n+1}^{N+n} \lambda(Z^*)_i &= \underset{W \in \mathbb{S}^{N+n}}{\text{minimize}} \quad \langle Z^*, W \rangle \\ &\text{subject to} \quad 0 \preceq W \preceq I \\ &\quad \text{tr } W = N \end{aligned} \quad (1872a)$$

which has an optimal solution that is known in closed form (p.533, §4.5.1.1). This eigenvalue sum is zero when  $Z^*$  has rank  $n$  or less.

Foreknowledge of optimal  $Z^*$ , to make possible this search for  $W$ , implies iteration; *id est*, semidefinite program (818) is solved for  $Z^*$  initializing  $W = I$  or  $W = \mathbf{0}$ . Once found,  $Z^*$  becomes constant in semidefinite program (1872a) where a new normal direction  $W$  is found as its optimal solution. Then this cycle (818) (1872a) iterates until convergence.

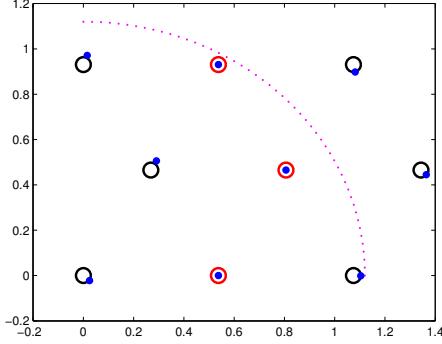


Figure 104: Typical solution for 3-lattice in Figure 99 with noise factor  $\eta = 0.1$  (821). Three red vertical middle nodes are anchors; remaining nodes are sensors. Radio range of sensor 1 indicated by arc; radius = 1.12. Actual sensor indicated by target  $\circ$  while its localization is indicated by bullet  $\bullet$ . Rank-2 solution found in 2 iterations (818) (1872a).

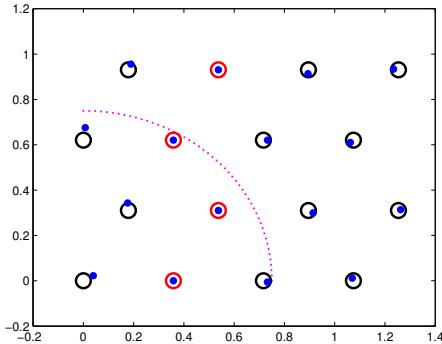


Figure 105: Typical solution for 4-lattice in Figure 100 with noise factor  $\eta = 0.1$  (821). Four red vertical middle-left nodes are anchors; remaining nodes are sensors. Radio range of sensor 1 indicated by arc; radius = 0.75. Actual sensor indicated by target  $\circ$  while its localization is indicated by bullet  $\bullet$ . Rank-2 solution found in 7 iterations (818) (1872a).

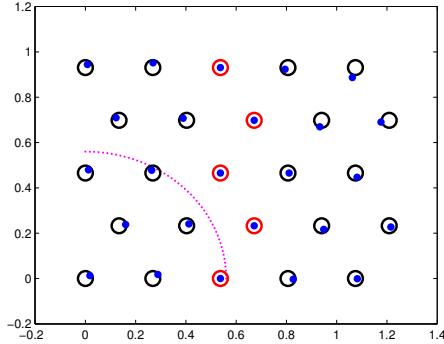


Figure 106: Typical solution for 5-lattice in Figure 101 with noise factor  $\eta = 0.1$  (821). Five red vertical middle nodes are anchors; remaining nodes are sensors. Radio range of sensor 1 indicated by arc; radius = 0.56. Actual sensor indicated by target  $\circ$  while its localization is indicated by bullet  $\bullet$ . Rank-2 solution found in 3 iterations (818) (1872a).

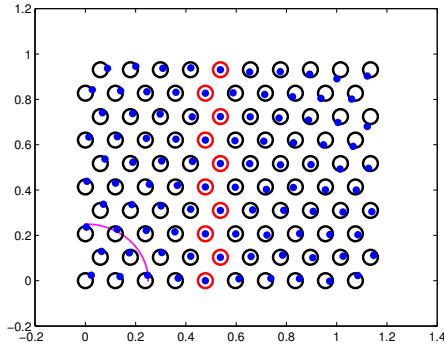


Figure 107: Typical solution for 10-lattice with noise factor  $\eta = 0.1$  (821) compares better than Carter & Jin [79, fig.4.2]. Ten red vertical middle nodes are anchors; the rest are sensors. Radio range of sensor 1 indicated by arc; radius = 0.25. Actual sensor indicated by target  $\circ$  while its localization is indicated by bullet  $\bullet$ . Rank-2 solution found in 5 iterations (818) (1872a).

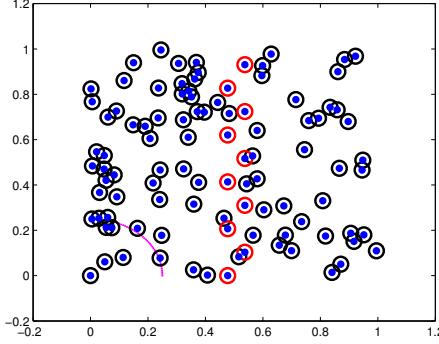


Figure 108: Typical localization of 100 randomized noiseless sensors ( $\eta = 0$  (821)) is exact despite incomplete EDM. Ten red vertical middle nodes are anchors; remaining nodes are sensors. Radio range of sensor at origin indicated by arc; radius = 0.25. Actual sensor indicated by target  $\circ$  while its localization is indicated by bullet  $\bullet$ . Rank-2 solution found in 3 iterations (818) (1872a).

When  $\text{rank } Z^* = n$ , solution via this convex iteration solves sensor-network localization problem (811) and its equivalent (816).

### numerical solution

In all examples to follow, number of anchors

$$m = \sqrt{N} \quad (820)$$

equals square root of cardinality  $N$  of list  $X$ . Indices set  $\mathcal{I}$  identifying all measurable distances  $\bullet$  is ascertained from connectivity matrix (812), (813), (814), or (815). We solve iteration (818) (1872a) in dimension  $n = 2$  for each respective example illustrated in Figure 98 through Figure 101.

In presence of negligible noise, true position is reliably localized for every standardized example; noteworthy insofar as each example represents an incomplete graph. This implies that the set of all optimal solutions having least rank must be small.

To make the examples interesting and consistent with previous work, we randomize each range of distance-square that bounds  $\langle G, (e_i - e_j)(e_i - e_j)^T \rangle$  in (818); *id est*, for each and every  $(i, j) \in \mathcal{I}$

$$\begin{aligned} \overline{d_{ij}} &= d_{ij}(1 + \sqrt{3}\eta\chi_l)^2 \\ \underline{d_{ij}} &= d_{ij}(1 - \sqrt{3}\eta\chi_{l+1})^2 \end{aligned} \quad (821)$$

where  $\eta = 0.1$  is a constant noise factor,  $\chi_l$  is the  $l^{\text{th}}$  sample of a noise process realization uniformly distributed in the interval  $(0, 1)$  like `rand(1)` from MATLAB, and  $d_{ij}$  is actual distance-square from  $i^{\text{th}}$  to  $j^{\text{th}}$  sensor. Because of distinct function calls to `rand()`, each range of distance-square  $[\underline{d_{ij}}, \overline{d_{ij}}]$  is not necessarily centered on actual distance-square  $d_{ij}$ . Unit stochastic variance is provided by factor  $\sqrt{3}$ .

Figure 103 through Figure 106 each illustrate one realization of numerical solution to the standardized lattice problems posed by Figure 98 through Figure 101 respectively. Exact localization, by any method, is impossible because of measurement noise. Certainly, by inspection of their published graphical data, our results are better than those of Carter & Jin. (Figure 107, 108, 109) Obviously our solutions do not suffer from those

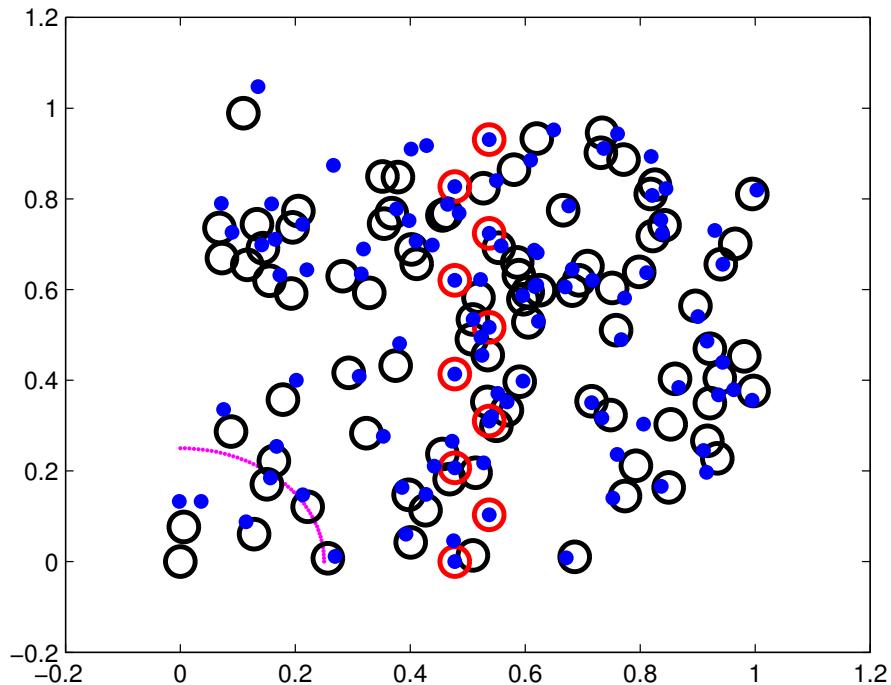


Figure 109: Typical solution for 100 randomized sensors with noise factor  $\eta = 0.1$  (821); worst measured average sensor error  $\approx 0.0044$  compares better than Carter & Jin's 0.0154 computed in 0.71s [79, p.19]. Ten red vertical middle nodes are anchors; same as before. Remaining nodes are sensors. Interior anchor placement makes localization difficult. Radio range of sensor at origin indicated by arc; radius = 0.25. Actual sensor indicated by target  $\circ$  while its localization is indicated by bullet  $\bullet$ . After 1 iteration rank  $G = 92$ , after 2 iterations rank  $G = 4$ . Rank-2 solution found in 3 iterations (818) (1872a). (Regular lattice in Figure 107 is actually harder to solve, requiring more iterations.) Runtime for SDPT3 [389] under cvx [191] is a few minutes on 2009 vintage laptop Core 2 Duo CPU (Intel T6400@2GHz, 800MHz FSB).

compaction-type errors (clustering of localized sensors) exhibited by Biswas' graphical results for the same noise factor  $\eta$ .

### localization example conclusion

Solution to this sensor-network localization problem became apparent by understanding geometry of optimization. Trace of a matrix, to a student of linear algebra, is perhaps a sum of eigenvalues. But to us, trace represents the normal  $I$  to some hyperplane in Euclidean vector space. (Figure 96)

Our solutions are globally optimal, requiring: 1) no centralized-gradient postprocessing heuristic refinement as in [47] because there is effectively no relaxation of (816) at global optimality, 2) no implicit postprojection on rank-2 positive semidefinite matrices induced by nonzero  $G - X^T X$  denoting suboptimality as occurs in [48] [49] [50] [79] [243] [251]; indeed,  $G^* = X^{*\top} X^*$  by convex iteration.

Numerical solution to noisy problems, containing sensor variables well in excess of 100, becomes difficult via the holistic semidefinite program we proposed. When problem size is within reach of contemporary general-purpose semidefinite program solvers, then the convex iteration we presented inherently overcomes limitations of Carter & Jin with respect to both noise performance and ability to localize in any desired affine dimension.

The legacy of Carter, Jin, Saunders, & Ye [79] is a sobering demonstration of the need for more efficient methods for solution of semidefinite programs, while that of So & Ye [357] forever bonds distance geometry to semidefinite programming. Elegance of our semidefinite problem statement (818), for constraining affine dimension of sensor-network localization, should provide some *impetus* to focus more research on computational intensity of general-purpose semidefinite program solvers. An approach different from interior-point methods is required; higher speed and greater accuracy from a simplex-like solver is what is needed.  $\square$

#### 4.5.1.2.5 Example. Nonnegative spectral factorization. (confer §3.14.2.0.2)

Having found optimal real coefficient vectors  $v^*, u^*$  for a sixteenth order magnitude square transfer function, evaluated along the  $j\omega$  axis (p.209),

$$|H(j\omega)|^2 = H(j\omega)H(-j\omega) = \frac{1 + v_1^* \omega^2 + v_2^* \omega^4 + \dots + v_8^* \omega^{16}}{1 + u_1^* \omega^2 + u_2^* \omega^4 + \dots + u_8^* \omega^{16}} \quad (666)$$

we wish to find real coefficients  $b, a$  for corresponding Fourier transform

$$H(j\omega) = \frac{1 + b_1 j\omega + b_2 (j\omega)^2 + \dots + b_8 (j\omega)^8}{1 + a_1 j\omega + a_2 (j\omega)^2 + \dots + a_8 (j\omega)^8} \quad (663)$$

These coefficients  $b, a, v^*, u^*$  are related through simultaneous nonlinear algebraic equations:

$$\begin{aligned} v_1^* &= b_1^2 - 2b_2, & u_1^* &= a_1^2 - 2a_2 \\ v_2^* &= b_2^2 - 2b_1b_3 + 2b_4, & u_2^* &= a_2^2 - 2a_1a_3 + 2a_4 \\ v_3^* &= b_3^2 - 2b_2b_4 + 2b_1b_5 - 2b_6, & u_3^* &= a_3^2 - 2a_2a_4 + 2a_1a_5 - 2a_6 \\ v_4^* &= b_4^2 - 2b_3b_5 + 2b_2b_6 - 2b_1b_7 + 2b_8, & u_4^* &= a_4^2 - 2a_3a_5 + 2a_2a_6 - 2a_1a_7 + 2a_8 \\ v_5^* &= b_5^2 - 2b_4b_6 + 2b_3b_7 - 2b_2b_8, & u_5^* &= a_5^2 - 2a_4a_6 + 2a_3a_7 - 2a_2a_8 \\ v_6^* &= b_6^2 - 2b_5b_7 + 2b_4b_8, & u_6^* &= a_6^2 - 2a_5a_7 + 2a_4a_8 \\ v_7^* &= b_7^2 - 2b_6b_8, & u_7^* &= a_7^2 - 2a_6a_8 \\ v_8^* &= b_8^2, & u_8^* &= a_8^2 \end{aligned} \quad (822)$$

Define a rank-one matrix

$$G(b) \triangleq \begin{bmatrix} 1 \\ b \end{bmatrix} [1 \ b^T] = \begin{bmatrix} 1 & b_1 & b_2 & b_3 & b_4 & b_5 & b_6 & b_7 & b_8 \\ b_1 & b_1^2 & b_1 b_2 & b_1 b_3 & b_1 b_4 & b_1 b_5 & b_1 b_6 & b_1 b_7 & b_1 b_8 \\ b_2 & b_1 b_2 & b_2^2 & b_2 b_3 & b_2 b_4 & b_2 b_5 & b_2 b_6 & b_2 b_7 & b_2 b_8 \\ b_3 & b_1 b_3 & b_2 b_3 & b_3^2 & b_3 b_4 & b_3 b_5 & b_3 b_6 & b_3 b_7 & b_3 b_8 \\ b_4 & b_1 b_4 & b_2 b_4 & b_3 b_4 & b_4^2 & b_4 b_5 & b_4 b_6 & b_4 b_7 & b_4 b_8 \\ b_5 & b_1 b_5 & b_2 b_5 & b_3 b_5 & b_4 b_5 & b_5^2 & b_5 b_6 & b_5 b_7 & b_5 b_8 \\ b_6 & b_1 b_6 & b_2 b_6 & b_3 b_6 & b_4 b_6 & b_5 b_6 & b_6^2 & b_6 b_7 & b_6 b_8 \\ b_7 & b_1 b_7 & b_2 b_7 & b_3 b_7 & b_4 b_7 & b_5 b_7 & b_6 b_7 & b_7^2 & b_7 b_8 \\ b_8 & b_1 b_8 & b_2 b_8 & b_3 b_8 & b_4 b_8 & b_5 b_8 & b_6 b_8 & b_7 b_8 & b_8^2 \end{bmatrix} \in \mathbb{S}^9 \quad (823)$$

(Matrix  $G(a)$  is similarly defined.) Observe that  $v^*$  in (822) is formed by summing antidiagonals of  $G(b)$  whose entries alternate sign. A particular sum is specified by a predetermined symmetric matrix constant  $A_i$  (confer (58)) from a set  $\{A_i \in \mathbb{S}^9, i=1 \dots 8\}$ . With

$$A = \begin{bmatrix} \text{svec}(A_1)^T \\ \vdots \\ \text{svec}(A_8)^T \end{bmatrix} \in \mathbb{R}^{8 \times 9(9+1)/2} \quad (698)$$

as previously defined in §4.1.1, all the sums (822) may be stated as two linear equalities  $A \text{svec } G(b) = v^*$  and  $A \text{svec } G(a) = u^*$ . Then the problem of finding coefficients  $b$  may be stated as a feasibility problem<sup>4.35</sup>

$$\begin{array}{ll} \text{find} & b \in \mathbb{R}^8 \\ \text{subject to} & A \text{svec } G = v^* \\ & \begin{bmatrix} 1 \\ b \end{bmatrix} = G(:, 1) \\ & b \succeq 0 \\ & (G \succeq 0) \\ & \text{rank } G = 1 \end{array} \quad (824)$$

The rank-one constraint is handled by convex iteration, as explained in §4.5.1. Positive semidefiniteness is parenthetical here because, for rank-one matrices, symmetry is necessary and sufficient (§A.3.1.0.7).  $\square$

#### 4.5.1.2.6 Example. Nonnegative spectral factorization II.

The purpose of spectral factorization, in electronics, is to facilitate high order filter implementation in the form of passive and active circuitry. Cascades of second-order (Laplace) sections are preferred because component sensitivity becomes manageable and because needed complex poles and zeros cannot be obtained from a first-order section.

Nonnegative spectral factorization on a magnitude square transfer function, evaluated along the  $j\omega$  axis, was performed in Example 4.5.1.2.5 to recover its corresponding Fourier transform.<sup>4.36</sup> In this example, we nonnegatively decompose a high order magnitude square transfer function into a product of successively lower order magnitude square transfer functions. Once fourth order magnitude square functions are found, then

<sup>4.35</sup> separately from the similar optimization problem to find vector  $a$ . Stability requires  $a \succeq 0$  with more constraints on  $a$ . Minimum phase requires  $b \succeq 0$  and more constraints on  $b$  that are missing from problem statement (824). Both stability and minimum phase may be enforced, subsequent to spectral factorization, by negating positive real parts of poles and zeros respectively in order to move them into the left half (Laplace)  $s$ -plane with no impact to  $|H(j\omega)|$ .

<sup>4.36</sup> When there are no poles on the  $j\omega$  axis, a Laplace transform can be recovered from a Fourier transform by substitution  $j\omega \leftarrow s$ .

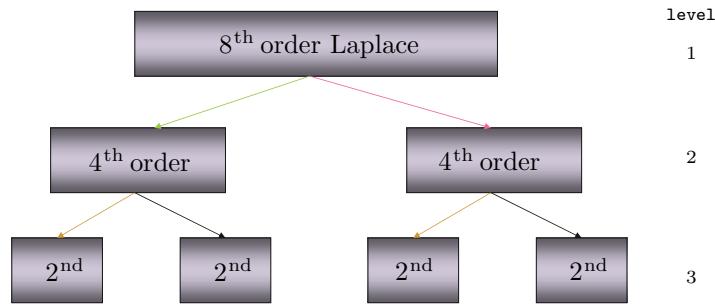


Figure 110: Nonnegative spectral factorization, high order bisection strategy.  $\eta = 8^{\text{th}}$  order Laplace transform corresponds to  $2\eta = 16^{\text{th}}$  order magnitude square transfer function. Because numerator  $v$  and denominator  $u$  are factored separately, number of factorizations  $= 2(\log_2(\eta) - 1)$ . In the text, double dots  $\ddot{v}, \ddot{u}$  connote first bifurcation (level 2). Triple dots  $\ddot{\ddot{v}}, \ddot{\ddot{u}}$  connote second bifurcations (level 3). Factors per level  $= 2^{\text{level}-1}$ .

corresponding second-order Laplace transfer function coefficients are ascertained from (665) and then passive component values can be determined from those coefficients.

Our strategy, for an eighth order Laplace transfer function, is illustrated in Figure 110. We begin at the tree's level 2 factorization. Nonnegative decomposition of a 16<sup>th</sup> order magnitude square transfer function into two 8<sup>th</sup> order functions

$$\frac{1 + v_1^* \omega^2 + v_2^* \omega^4 + \dots + v_8^* \omega^{16}}{1 + u_1^* \omega^2 + u_2^* \omega^4 + \dots + u_8^* \omega^{16}} = \frac{1 + \ddot{v}_1 \omega^2 + \ddot{v}_2 \omega^4 + \ddot{v}_3 \omega^6 + \ddot{v}_4 \omega^8}{1 + \ddot{u}_1 \omega^2 + \ddot{u}_2 \omega^4 + \ddot{u}_3 \omega^6 + \ddot{u}_4 \omega^8} \frac{1 + \ddot{v}_5 \omega^2 + \ddot{v}_6 \omega^4 + \ddot{v}_7 \omega^6 + \ddot{v}_8 \omega^8}{1 + \ddot{u}_5 \omega^2 + \ddot{u}_6 \omega^4 + \ddot{u}_7 \omega^6 + \ddot{u}_8 \omega^8} \quad (825)$$

implies these simultaneous algebraic identifications with known real coefficient vectors  $v^*, u^*$ :

$$\begin{aligned}
v_1^* &= \ddot{v}_1 + \ddot{v}_5, & u_1^* &= \ddot{u}_1 + \ddot{u}_5 \\
v_2^* &= \ddot{v}_2 + \ddot{v}_6 + \ddot{v}_1\ddot{v}_5, & u_2^* &= \ddot{u}_2 + \ddot{u}_6 + \ddot{u}_1\ddot{u}_5 \\
v_3^* &= \ddot{v}_3 + \ddot{v}_7 + \ddot{v}_1\ddot{v}_6 + \ddot{v}_2\ddot{v}_5, & u_3^* &= \ddot{u}_3 + \ddot{u}_7 + \ddot{u}_1\ddot{u}_6 + \ddot{u}_2\ddot{u}_5 \\
v_4^* &= \ddot{v}_4 + \ddot{v}_8 + \ddot{v}_1\ddot{v}_7 + \ddot{v}_2\ddot{v}_6 + \ddot{v}_3\ddot{v}_5, & u_4^* &= \ddot{u}_4 + \ddot{u}_8 + \ddot{u}_1\ddot{u}_7 + \ddot{u}_2\ddot{u}_6 + \ddot{u}_3\ddot{u}_5 \\
v_5^* &= \ddot{v}_4\ddot{v}_5 + \ddot{v}_3\ddot{v}_6 + \ddot{v}_2\ddot{v}_7 + \ddot{v}_1\ddot{v}_8, & u_5^* &= \ddot{u}_4\ddot{v}_5 + \ddot{u}_3\ddot{v}_6 + \ddot{u}_2\ddot{v}_7 + \ddot{u}_1\ddot{v}_8 \\
v_6^* &= \ddot{v}_4\ddot{v}_6 + \ddot{v}_3\ddot{v}_7 + \ddot{v}_2\ddot{v}_8, & u_6^* &= \ddot{u}_4\ddot{v}_6 + \ddot{u}_3\ddot{v}_7 + \ddot{u}_2\ddot{v}_8 \\
v_7^* &= \ddot{v}_4\ddot{v}_7 + \ddot{v}_3\ddot{v}_8, & u_7^* &= \ddot{u}_4\ddot{v}_7 + \ddot{u}_3\ddot{v}_8 \\
v_8^* &= \ddot{v}_4\ddot{v}_8, & u_8^* &= \ddot{u}_4\ddot{v}_8
\end{aligned} \tag{826}$$

Now define a rank-one matrix for the numerator

$$G(\ddot{v}) \triangleq \begin{bmatrix} 1 \\ \ddot{v} \end{bmatrix} \begin{bmatrix} 1 & \ddot{v}^T \end{bmatrix} = \begin{bmatrix} 1 & \ddot{v}_1 & \ddot{v}_2 & \ddot{v}_3 & \ddot{v}_4 & \ddot{v}_5 & \ddot{v}_6 & \ddot{v}_7 & \ddot{v}_8 \\ \ddot{v}_1 & \ddot{v}_1^2 & \ddot{v}_1\ddot{v}_2 & \ddot{v}_1\ddot{v}_3 & \ddot{v}_1\ddot{v}_4 & \ddot{v}_1\ddot{v}_5 & \ddot{v}_1\ddot{v}_6 & \ddot{v}_1\ddot{v}_7 & \ddot{v}_1\ddot{v}_8 \\ \ddot{v}_2 & \ddot{v}_1\ddot{v}_2 & \ddot{v}_2^2 & \ddot{v}_2\ddot{v}_3 & \ddot{v}_2\ddot{v}_4 & \ddot{v}_2\ddot{v}_5 & \ddot{v}_2\ddot{v}_6 & \ddot{v}_2\ddot{v}_7 & \ddot{v}_2\ddot{v}_8 \\ \ddot{v}_3 & \ddot{v}_1\ddot{v}_3 & \ddot{v}_2\ddot{v}_3 & \ddot{v}_3^2 & \ddot{v}_3\ddot{v}_4 & \ddot{v}_3\ddot{v}_5 & \ddot{v}_3\ddot{v}_6 & \ddot{v}_3\ddot{v}_7 & \ddot{v}_3\ddot{v}_8 \\ \ddot{v}_4 & \ddot{v}_1\ddot{v}_4 & \ddot{v}_2\ddot{v}_4 & \ddot{v}_3\ddot{v}_4 & \ddot{v}_4^2 & \ddot{v}_4\ddot{v}_5 & \ddot{v}_4\ddot{v}_6 & \ddot{v}_4\ddot{v}_7 & \ddot{v}_4\ddot{v}_8 \\ \ddot{v}_5 & \ddot{v}_1\ddot{v}_5 & \ddot{v}_2\ddot{v}_5 & \ddot{v}_3\ddot{v}_5 & \ddot{v}_4\ddot{v}_5 & \ddot{v}_5^2 & \ddot{v}_5\ddot{v}_6 & \ddot{v}_5\ddot{v}_7 & \ddot{v}_5\ddot{v}_8 \\ \ddot{v}_6 & \ddot{v}_1\ddot{v}_6 & \ddot{v}_2\ddot{v}_6 & \ddot{v}_3\ddot{v}_6 & \ddot{v}_4\ddot{v}_6 & \ddot{v}_5\ddot{v}_6 & \ddot{v}_6^2 & \ddot{v}_6\ddot{v}_7 & \ddot{v}_6\ddot{v}_8 \\ \ddot{v}_7 & \ddot{v}_1\ddot{v}_7 & \ddot{v}_2\ddot{v}_7 & \ddot{v}_3\ddot{v}_7 & \ddot{v}_4\ddot{v}_7 & \ddot{v}_5\ddot{v}_7 & \ddot{v}_6\ddot{v}_7 & \ddot{v}_7^2 & \ddot{v}_7\ddot{v}_8 \\ \ddot{v}_8 & \ddot{v}_1\ddot{v}_8 & \ddot{v}_2\ddot{v}_8 & \ddot{v}_3\ddot{v}_8 & \ddot{v}_4\ddot{v}_8 & \ddot{v}_5\ddot{v}_8 & \ddot{v}_6\ddot{v}_8 & \ddot{v}_7\ddot{v}_8 & \ddot{v}_8^2 \end{bmatrix} \in \mathbb{S}^9 \quad (827)$$

(Matrix  $G(\ddot{u})$  is defined similarly for the denominator.) Terms in (826) are picked out of  $G(\ddot{v})$  by a predetermined symmetric matrix constant  $\ddot{A}_i$  (confer (58)) from a set

$\{\ddot{A}_i \in \mathbb{S}^9, i=1 \dots 8\}$ . Populating rows of

$$A = \begin{bmatrix} \text{svec}(\ddot{A}_1)^T \\ \vdots \\ \text{svec}(\ddot{A}_8)^T \end{bmatrix} \in \mathbb{R}^{8 \times 9(9+1)/2} \quad (698)$$

with vectorized  $\ddot{A}_i$  (as in §4.1.1), sums (826) are succinctly represented by two linear equalities  $A \text{svec } G(\ddot{v}) = v^*$  and  $A \text{svec } G(\ddot{u}) = u^*$ . Then this spectral factorization in  $\ddot{v}$  may be posed as a feasibility problem

$$\begin{array}{ll} \underset{G \in \mathbb{S}^9}{\text{find}} & \ddot{v} \in \mathbb{R}^8 \\ \text{subject to} & A \text{svec } G = v^* \\ & \begin{bmatrix} 1 \\ \ddot{v} \end{bmatrix} = G(:, 1) \\ & \ddot{v} \succeq 0 \\ & (G \succeq 0) \\ & \text{rank } G = 1 \end{array} \quad (828)$$

Having found two 8<sup>th</sup> order square spectral factors in nonnegative  $\ddot{v}^*$  from (828), two pairs of 4<sup>th</sup> order **level 3** factors remain to be found:

$$\frac{1 + \ddot{v}_1^* \omega^2 + \ddot{v}_2^* \omega^4 + \ddot{v}_3^* \omega^6 + \ddot{v}_4^* \omega^8}{1 + \ddot{u}_1^* \omega^2 + \ddot{u}_2^* \omega^4 + \ddot{u}_3^* \omega^6 + \ddot{u}_4^* \omega^8} = \frac{1 + \ddot{v}_1 \omega^2 + \ddot{v}_2 \omega^4}{1 + \ddot{u}_1 \omega^2 + \ddot{u}_2 \omega^4} \frac{1 + \ddot{v}_3 \omega^2 + \ddot{v}_4 \omega^4}{1 + \ddot{u}_3 \omega^2 + \ddot{u}_4 \omega^4} \quad (829)$$

$$\frac{1 + \ddot{v}_5^* \omega^2 + \ddot{v}_6^* \omega^4 + \ddot{v}_7^* \omega^6 + \ddot{v}_8^* \omega^8}{1 + \ddot{u}_5^* \omega^2 + \ddot{u}_6^* \omega^4 + \ddot{u}_7^* \omega^6 + \ddot{u}_8^* \omega^8} = \frac{1 + \ddot{v}_5 \omega^2 + \ddot{v}_6 \omega^4}{1 + \ddot{u}_5 \omega^2 + \ddot{u}_6 \omega^4} \frac{1 + \ddot{v}_7 \omega^2 + \ddot{v}_8 \omega^4}{1 + \ddot{u}_7 \omega^2 + \ddot{u}_8 \omega^4} \quad (830)$$

$$\begin{array}{ll} \ddot{v}_1^* = \ddot{v}_1 + \ddot{v}_3, & \ddot{u}_1^* = \ddot{u}_1 + \ddot{u}_3 \\ \ddot{v}_2^* = \ddot{v}_2 + \ddot{v}_4 + \ddot{v}_1 \ddot{v}_3, & \ddot{u}_2^* = \ddot{u}_2 + \ddot{u}_4 + \ddot{u}_1 \ddot{u}_3 \\ \ddot{v}_3^* = \ddot{v}_1 \ddot{v}_4 + \ddot{v}_2 \ddot{v}_3, & \ddot{u}_3^* = \ddot{u}_1 \ddot{v}_4 + \ddot{u}_2 \ddot{u}_3 \\ \ddot{v}_4^* = \ddot{v}_2 \ddot{v}_4, & \ddot{u}_4^* = \ddot{u}_2 \ddot{u}_4 \end{array} \quad (831)$$

$$\begin{array}{ll} \ddot{v}_5^* = \ddot{v}_5 + \ddot{v}_7, & \ddot{u}_5^* = \ddot{u}_5 + \ddot{u}_7 \\ \ddot{v}_6^* = \ddot{v}_6 + \ddot{v}_8 + \ddot{v}_5 \ddot{v}_7, & \ddot{u}_6^* = \ddot{u}_6 + \ddot{u}_8 + \ddot{u}_5 \ddot{u}_7 \\ \ddot{v}_7^* = \ddot{v}_5 \ddot{v}_8 + \ddot{v}_6 \ddot{v}_7, & \ddot{u}_7^* = \ddot{u}_5 \ddot{v}_8 + \ddot{u}_6 \ddot{u}_7 \\ \ddot{v}_8^* = \ddot{v}_6 \ddot{v}_8, & \ddot{u}_8^* = \ddot{u}_6 \ddot{u}_8 \end{array} \quad (832)$$

$$G(\ddot{v}) \triangleq \begin{bmatrix} 1 & \ddot{v}^T \end{bmatrix} \begin{bmatrix} 1 & \ddot{v}^T \end{bmatrix}^T = \begin{bmatrix} 1 & \ddot{v}_1 & \ddot{v}_2 & \ddot{v}_3 & \ddot{v}_4 & \ddot{v}_5 & \ddot{v}_6 & \ddot{v}_7 & \ddot{v}_8 \\ \ddot{v}_1 & \ddot{v}_1^2 & \ddot{v}_1 \ddot{v}_2 & \ddot{v}_1 \ddot{v}_3 & \ddot{v}_1 \ddot{v}_4 & \ddot{v}_1 \ddot{v}_5 & \ddot{v}_1 \ddot{v}_6 & \ddot{v}_1 \ddot{v}_7 & \ddot{v}_1 \ddot{v}_8 \\ \ddot{v}_2 & \ddot{v}_1 \ddot{v}_2 & \ddot{v}_2^2 & \ddot{v}_2 \ddot{v}_3 & \ddot{v}_2 \ddot{v}_4 & \ddot{v}_2 \ddot{v}_5 & \ddot{v}_2 \ddot{v}_6 & \ddot{v}_2 \ddot{v}_7 & \ddot{v}_2 \ddot{v}_8 \\ \ddot{v}_3 & \ddot{v}_1 \ddot{v}_3 & \ddot{v}_2 \ddot{v}_3 & \ddot{v}_3^2 & \ddot{v}_3 \ddot{v}_4 & \ddot{v}_3 \ddot{v}_5 & \ddot{v}_3 \ddot{v}_6 & \ddot{v}_3 \ddot{v}_7 & \ddot{v}_3 \ddot{v}_8 \\ \ddot{v}_4 & \ddot{v}_1 \ddot{v}_4 & \ddot{v}_2 \ddot{v}_4 & \ddot{v}_3 \ddot{v}_4 & \ddot{v}_4^2 & \ddot{v}_4 \ddot{v}_5 & \ddot{v}_4 \ddot{v}_6 & \ddot{v}_4 \ddot{v}_7 & \ddot{v}_4 \ddot{v}_8 \\ \ddot{v}_5 & \ddot{v}_1 \ddot{v}_5 & \ddot{v}_2 \ddot{v}_5 & \ddot{v}_3 \ddot{v}_5 & \ddot{v}_4 \ddot{v}_5 & \ddot{v}_5^2 & \ddot{v}_5 \ddot{v}_6 & \ddot{v}_5 \ddot{v}_7 & \ddot{v}_5 \ddot{v}_8 \\ \ddot{v}_6 & \ddot{v}_1 \ddot{v}_6 & \ddot{v}_2 \ddot{v}_6 & \ddot{v}_3 \ddot{v}_6 & \ddot{v}_4 \ddot{v}_6 & \ddot{v}_5 \ddot{v}_6 & \ddot{v}_6^2 & \ddot{v}_6 \ddot{v}_7 & \ddot{v}_6 \ddot{v}_8 \\ \ddot{v}_7 & \ddot{v}_1 \ddot{v}_7 & \ddot{v}_2 \ddot{v}_7 & \ddot{v}_3 \ddot{v}_7 & \ddot{v}_4 \ddot{v}_7 & \ddot{v}_5 \ddot{v}_7 & \ddot{v}_6 \ddot{v}_7 & \ddot{v}_7^2 & \ddot{v}_7 \ddot{v}_8 \\ \ddot{v}_8 & \ddot{v}_1 \ddot{v}_8 & \ddot{v}_2 \ddot{v}_8 & \ddot{v}_3 \ddot{v}_8 & \ddot{v}_4 \ddot{v}_8 & \ddot{v}_5 \ddot{v}_8 & \ddot{v}_6 \ddot{v}_8 & \ddot{v}_7 \ddot{v}_8 & \ddot{v}_8^2 \end{bmatrix} \in \mathbb{S}^9 \quad (833)$$

Setting

$$A = \begin{bmatrix} \text{svec}(\ddot{A}_1)^T \\ \vdots \\ \text{svec}(\ddot{A}_8)^T \end{bmatrix} \in \mathbb{R}^{8 \times 9(9+1)/2} \quad (698)$$

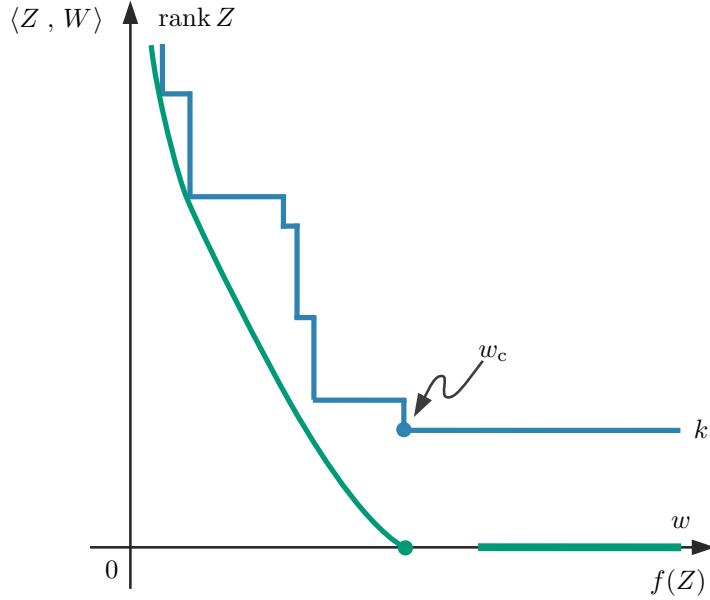


Figure 111: Regularization curve, parametrized by weight  $w$  for real convex objective  $f$  minimization (835) with rank constraint to  $k$  by convex iteration, illustrates discontinuity in  $f$ .

then all level 3 (Figure 110) nonnegative spectral factorization coefficients  $\ddot{v}$  are found at once by solving

$$\begin{aligned}
 & \underset{G \in \mathbb{S}^9}{\text{find}} \quad \ddot{v} \in \mathbb{R}^8 \\
 & \text{subject to} \quad A \text{ svec } G = \ddot{v}^* \\
 & \qquad \left[ \begin{array}{c} 1 \\ \ddot{v} \end{array} \right] = G(:, 1) \\
 & \qquad \ddot{v} \succeq 0 \\
 & \qquad (G \succeq 0) \\
 & \qquad \text{rank } G = 1
 \end{aligned} \tag{834}$$

The feasibility problem to find  $\ddot{u}$  is similar. All second-order Laplace transfer function coefficients can be found via (665).  $\square$

#### 4.5.2 regularization

We test the convex iteration technique, for constraining rank, over a wide range of problems beyond localization of randomized positions (Figure 109); e.g., stress (§7.2.2.7.1), ball packing (§5.4.2.2.6), and cardinality (§4.7). We have had some success introducing the direction matrix inner-product (819) as a regularization term<sup>4.37</sup>

$$\begin{aligned}
 & \underset{Z \in \mathbb{S}^N}{\text{minimize}} \quad f(Z) + w \langle Z, W \rangle \\
 & \text{subject to} \quad Z \in \mathcal{C} \\
 & \qquad Z \succeq 0
 \end{aligned} \tag{835}$$

<sup>4.37</sup> called *multiobjective-* or *vector optimization*. Proof of convergence for this convex iteration is identical to that in §4.5.1.2.1 because  $f$  is a convex real function, hence bounded below, and  $f(Z^*)$  is constant in (836).

$$\begin{aligned} & \underset{W \in \mathbb{S}^N}{\text{minimize}} \quad f(Z^*) + w \langle Z^*, W \rangle \\ & \text{subject to} \quad 0 \preceq W \preceq I \\ & \quad \text{tr } W = N - n \end{aligned} \tag{836}$$

whose purpose is to constrain rank, affine dimension, or cardinality:

The abstraction, that is Figure 111, is a synopsis; a broad generalization of accumulated empirical evidence: There exists a critical (smallest) weight  $w_c$  for which a rank constraint is just met. Graphical discontinuity can subsequently exist when there is a range of greater  $w$  providing required rank  $k$  but not necessarily increasing a minimization objective function  $f$ ; e.g., §4.7.0.0.2. Positive scalar  $w$  is chosen via bisection so that  $\langle Z^*, W^* \rangle$  just vanishes.

## 4.6 Constraining cardinality

The convex iteration technique for constraining rank can be applied to cardinality problems. There are parallels in its development analogous to how prototypical semidefinite program (697) resembles linear program (696) on page 222 [444]:

### 4.6.1 nonnegative variable

Our goal has been to reliably constrain rank in a semidefinite program. There is a direct analogy to linear programming that is simpler to present but, perhaps, more difficult to solve. In Optimization, that analogy is known as the *cardinality problem*.

Consider a feasibility problem  $Ax = b$ , but with an upper bound  $k$  on cardinality  $\|x\|_0$  of a nonnegative solution  $x$ : for  $A \in \mathbb{R}^{m \times n}$  and vector  $b \in \mathcal{R}(A)$

$$\begin{aligned} & \text{find} \quad x \in \mathbb{R}^n \\ & \text{subject to} \quad Ax = b \\ & \quad x \succeq 0 \\ & \quad \|x\|_0 \leq k \end{aligned} \tag{535}$$

where  $\|x\|_0 \leq k$  means<sup>4.38</sup> vector  $x$  has at most  $k$  nonzero entries; such a vector is presumed existent in the feasible set. Nonnegativity constraint  $x \succeq 0$  is analogous to positive semidefiniteness; the notation means vector  $x$  belongs to the nonnegative orthant  $\mathbb{R}_+^n$ . Cardinality is quasiconcave on  $\mathbb{R}_+^n$  just as rank is quasiconcave on  $\mathbb{S}_+^n$ . [65, §3.4.2]

#### 4.6.1.1 direction vector

We propose that cardinality-constrained feasibility problem (535) can be equivalently expressed as iteration of a sequence of two convex problems: for  $0 \leq k \leq n-1$

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \langle x, y \rangle \\ & \text{subject to} \quad Ax = b \\ & \quad x \succeq 0 \end{aligned} \tag{158}$$

$$\begin{aligned} \sum_{i=k+1}^n \pi(x^*)_i &= \underset{y \in \mathbb{R}^n}{\text{minimize}} \quad \langle x^*, y \rangle \\ \text{subject to} \quad 0 \preceq y \preceq \mathbf{1} \\ & \quad y^T \mathbf{1} = n - k \end{aligned} \tag{530}$$

---

<sup>4.38</sup>Although it is a metric (§5.2), cardinality  $\|x\|_0$  cannot be a norm (§3.2) because it is not positively homogeneous.

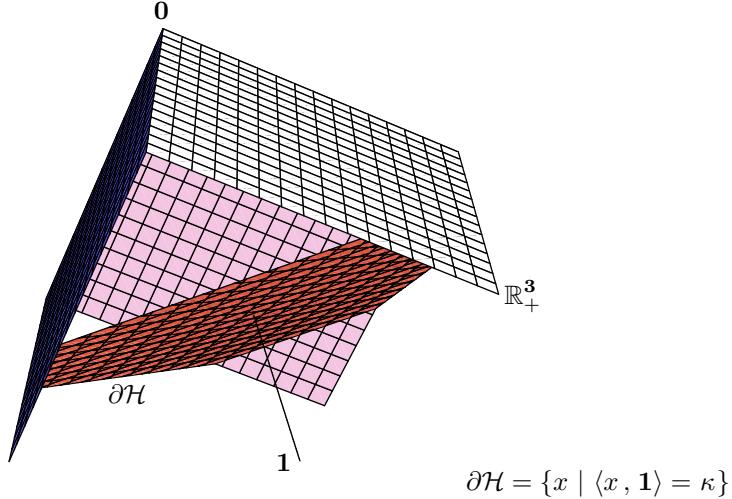


Figure 112: (confer Figure 96) 1-norm heuristic for cardinality minimization can be interpreted as minimization of a hyperplane,  $\partial\mathcal{H}$  with normal  $\mathbf{1}$ , over nonnegative orthant drawn here in  $\mathbb{R}^3$ . Polar of direction vector  $y = \mathbf{1}$  points toward origin.

where  $\pi$  is the (nonincreasing) *presorting function* (1471). This sequence is iterated until  $x^{*\top}y^*$  vanishes; *id est*, until desired cardinality is achieved. But this *global optimality* cannot be guaranteed.<sup>4.39</sup>

Problem (530) is analogous to the rank constraint problem; (p.248)

$$\begin{aligned} \sum_{i=k+1}^N \lambda(G^*)_i &= \underset{W \in \mathbb{S}^N}{\text{minimize}} \quad \langle G^*, W \rangle \\ \text{subject to} \quad 0 \preceq W \preceq I \\ \text{tr } W = N - k \end{aligned} \tag{1872a}$$

The feasible set of (530) is Linear Program's analogue to Fantope (§2.3.2.0.1); its optimal subset comprises a sum of  $n-k$  smallest entries from vector  $x$ . In context of problem (535), we want  $n-k$  entries of  $x$  to sum to zero; *id est*, we want a globally optimal objective  $x^{*\top}y^*$  to vanish: more generally, (confer(802))

$$\sum_{i=k+1}^n \pi(|x^*|)_i = \langle |x^*|, y^* \rangle = |x^*|^T y^* \triangleq 0 \tag{837}$$

defines *global optimality* for the iteration. Then  $n-k$  entries of  $x^*$  are themselves zero whenever their absolute sum is, and cardinality of  $x^* \in \mathbb{R}^n$  is at most  $k$ . *Optimal direction vector*  $y^*$  is defined as any nonnegative vector for which

$$\begin{aligned} \text{find} \quad x \in \mathbb{R}^n \\ \text{subject to} \quad Ax = b \\ x \succeq 0 \\ \|x\|_0 \leq k \end{aligned} \quad \equiv \quad \begin{aligned} \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \langle x, y^* \rangle \\ \text{subject to} \quad Ax = b \\ x \succeq 0 \end{aligned} \tag{158}$$

Existence of such a  $y^*$ , whose nonzero entries are complementary to those of  $x^*$ , is obvious assuming existence of a cardinality- $k$  solution  $x^*$ .

<sup>4.39</sup>When it succeeds, a sequence may be regarded as a *homotopy* to minimal 0-norm.

#### 4.6.1.2 direction vector interpretation

(confer §4.5.1.1) Vector  $y$  may be interpreted as a negative search direction; it points opposite to direction of movement of hyperplane  $\{x \mid \langle x, y \rangle = \tau\}$  in a minimization of real linear function  $\langle x, y \rangle$  over the feasible set in linear program (158). (p.62) Direction vector  $y$  is not unique. The feasible set of direction vectors in (530) is the convex hull of all cardinality- $(n-k)$  one-vectors; *videlicet*,

$$\text{conv}\{u \in \mathbb{R}^n \mid \text{card } u = n - k, u_i \in \{0, 1\}\} = \{a \in \mathbb{R}^n \mid \mathbf{1} \succeq a \succeq 0, \langle \mathbf{1}, a \rangle = n - k\} \quad (838)$$

This set, argument to  $\text{conv}\{\}$ , comprises the extreme points of set (838) which is a *nonnegative hypercube slice*. An optimal solution  $y$  to (530), that is an extreme point of its feasible set, is known in closed form: it has 1 in each entry corresponding to the  $n-k$  smallest entries of  $x^*$  and has 0 elsewhere. That particular polar direction  $-y$  can be interpreted<sup>4.40</sup> (by Proposition 7.1.3.0.3) as pointing toward the nonnegative orthant in the *Cartesian subspace*, whose basis is a subset of the Cartesian axes, containing all cardinality  $k$  (or less) vectors having the same ordering as  $x^*$ . Consequently, for that closed-form solution, (confer(803))

$$\sum_{i=k+1}^n \pi(|x^*|)_i = \langle |x^*|, y \rangle = |x^*|^T y \geq 0 \quad (839)$$

When  $y = \mathbf{1}$ , as in 1-norm minimization for example, then polar direction  $-y$  points directly at the origin (the cardinality-0 nonnegative vector) as in Figure 112. We sometimes solve (530) instead of employing a known closed form because a direction vector is not unique. Setting direction vector  $y$  instead in accordance with an iterative inverse weighting scheme, called *reweighting* [184], was described for the 1-norm by Huo [234, §4.11.3] in 1999.

#### 4.6.1.3 convergence can mean stalling

Convex iteration (158) (530) always converges to a *locally optimal solution*, a fixed point of possibly infeasible cardinality, by virtue of a monotonically nonincreasing real objective sequence. [289, §1.2] [43, §1.1] There can be no proof of global optimality, defined by (837). Constraining cardinality (solution to problem (535)) can often be achieved, but simple examples can be contrived that *stall* at a fixed point of infeasible cardinality; at a positive objective value  $\langle x^*, y \rangle = \tau > 0$ . Direction vector  $y$  is then manipulated, as countermeasure, to steer out of local minima; *e.g.*, complete randomization as in Example 4.6.1.5.1, or reinitialization to a random cardinality- $(n-k)$  vector in the same nonnegative orthant face demanded by the current iterate:  $y$  has nonnegative uniformly distributed random entries in  $(0, 1]$  corresponding to the  $n-k$  smallest entries of  $x^*$  and has 0 elsewhere. Zero entries behave like memory or state while randomness greatly diminishes likelihood of a stall. When this particular heuristic is successful, cardinality and objective sequence  $\langle x^*, y \rangle$  *versus* iteration are characterized by noisy monotonicity.

---

<sup>4.40</sup>Convex iteration (158) (530) is not a projection method because there is no thresholding or discard of variable-vector  $x$  entries. An optimal direction vector  $y$  must always reside on the feasible set boundary in (530) page 271; *id est*, it is ill-advised to attempt simultaneous optimization of variables  $x$  and  $y$ .

#### 4.6.1.4 algebraic derivation of direction vector for convex iteration

In §3.2.2.1.3, the compressed sensing problem was precisely represented as a nonconvex difference of convex functions bounded below by 0

$$\begin{array}{ll} \text{find} & x \in \mathbb{R}^n \\ \text{subject to} & \begin{array}{l} Ax = b \\ x \succeq 0 \\ \|x\|_0 \leq k \end{array} \end{array} \equiv \begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \|x\|_1 - \|x\|_k \\ \text{subject to} & \begin{array}{l} Ax = b \\ x \succeq 0 \end{array} \end{array} \quad (535)$$

where convex  $k$ -largest norm  $\|x\|_k$  is monotonic on  $\mathbb{R}_+^n$ . There we showed how (535) is equivalently stated in terms of gradients

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \langle x, \nabla \|x\|_1 - \nabla \|x\|_k \rangle \\ \text{subject to} & \begin{array}{l} Ax = b \\ x \succeq 0 \end{array} \end{array} \quad (840)$$

because

$$\|x\|_1 = x^T \nabla \|x\|_1, \quad \|x\|_k = x^T \nabla \|x\|_k, \quad x \succeq 0 \quad (841)$$

The objective function from (840) is a directional derivative (at  $x$  in direction  $x$ , §D.1.6, confer §D.1.4.1.1) of the objective function from (535) while the direction vector of convex iteration

$$y = \nabla \|x\|_1 - \nabla \|x\|_k \quad (842)$$

is an objective gradient where  $\nabla \|x\|_1 = \nabla \mathbf{1}^T x = \mathbf{1}$  under nonnegativity and

$$\nabla \|x\|_k = \nabla z^T x = \left. \begin{array}{l} \arg \max_{z \in \mathbb{R}^n} z^T x \\ \text{subject to} \quad \begin{array}{l} 0 \preceq z \preceq \mathbf{1} \\ z^T \mathbf{1} = k \end{array} \end{array} \right\}, \quad x \succeq 0 \quad (538)$$

is not unique. Substituting  $\mathbf{1} - z \leftarrow z$  the direction vector becomes

$$\begin{array}{ll} y = \mathbf{1} - \arg \max_{z \in \mathbb{R}^n} z^T x & \leftarrow \arg \min_{z \in \mathbb{R}^n} z^T x \\ \text{subject to} \quad \begin{array}{l} 0 \preceq z \preceq \mathbf{1} \\ z^T \mathbf{1} = k \end{array} & \text{subject to} \quad \begin{array}{l} 0 \preceq z \preceq \mathbf{1} \\ z^T \mathbf{1} = n - k \end{array} \end{array} \quad (530)$$

#### 4.6.1.5 optimality conditions for minimal cardinality

Now we see how global optimality conditions can be stated without reference to a dual problem: From conditions (474) for optimality of (535), it is necessary [65, §5.5.3] that

$$\begin{array}{ll} x^* \succeq 0 & (1) \\ Ax^* = b & (2) \\ \nabla \|x^*\|_1 - \nabla \|x^*\|_k + A^T \nu^* \succeq 0 & (3) \\ \langle \nabla \|x^*\|_1 - \nabla \|x^*\|_k + A^T \nu^*, x^* \rangle = 0 & (4L) \end{array} \quad (843)$$

These conditions must hold at any optimal solution (locally or globally). By (841), the fourth condition is identical to

$$\|x^*\|_1 - \|x^*\|_k + \nu^{*T} Ax^* = 0 \quad (4L) \quad (844)$$

Because a 1-norm

$$\|x\|_1 = \|x\|_k + \|\pi(|x|)_{k+1:n}\|_1 \quad (845)$$

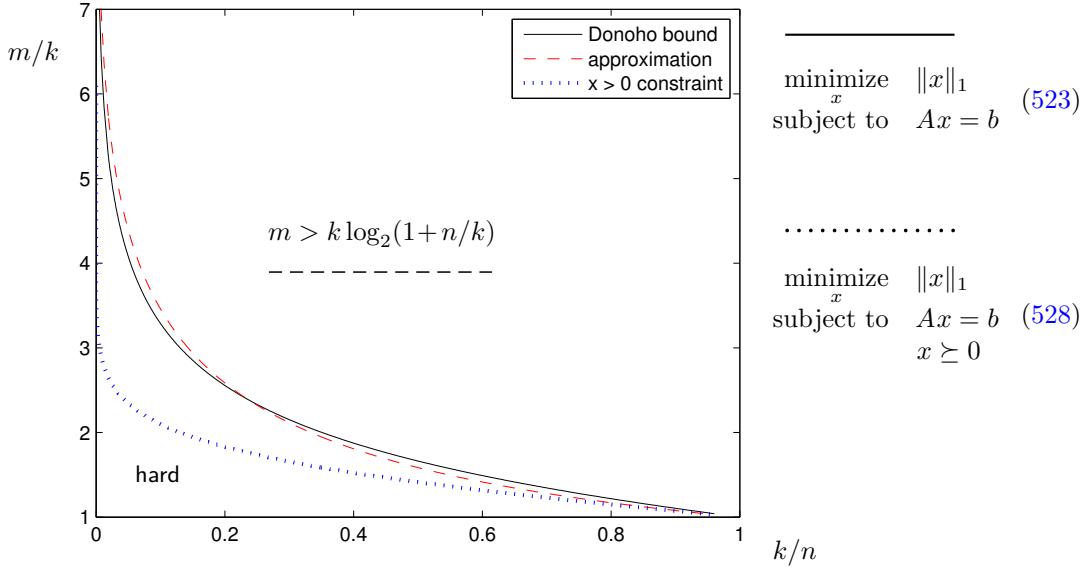


Figure 113: (confer Figure 75) For Gaussian random matrix  $A \in \mathbb{R}^{m \times n}$ , graph illustrates Donoho/Tanner least lower bound on number of measurements  $m$  below which recovery of  $k$ -sparse  $n$ -length signal  $x$  by linear programming fails with overwhelming probability. Hard problems are below curve, but not the reverse; *id est*, failure above depends on proximity. Inequality demarcates approximation (---) to empirical phase transition from [24]. Problems having nonnegativity constraint (···) are easier to solve. [140] [141]

is separable into  $k$  largest and  $n - k$  smallest absolute entries,

$$\|\pi(|x|)_{k+1:n}\|_1 = 0 \Leftrightarrow \|x\|_0 \leq k \quad (4g)$$

is a necessary condition for global optimality. By assumption, matrix  $A$  is wide and  $b \neq \mathbf{0} \Rightarrow Ax^* \neq \mathbf{0}$ . This means  $\nu^* \in \mathcal{N}(A^T) \subset \mathbb{R}^m$ , and  $\nu^* = \mathbf{0}$  when  $A$  is full-rank. By definition,  $\nabla\|x\|_1 \succeq \nabla\|x\|_k$  always holds. Assuming existence of a cardinality- $k$  solution, then only three of the four conditions are necessary and sufficient for global optimality of (535):

$$\begin{array}{ll} x^* \succeq 0 & (1) \\ Ax^* = b & (2) \\ \|x^*\|_1 - \|x^*\|_k = 0 & (4g) \end{array} \quad (847)$$

meaning, global optimality of a feasible solution to (535) is identified by a zero objective.

#### 4.6.1.5.1 Example. Sparsest solution to $Ax = b$ . [77] [136]

(confer Example 4.6.2.0.2) Data (728) induces sparsest solution not easily recoverable by least 1-norm; *id est*, not by compressed sensing because of proximity to a theoretical lower bound on number of *measurements*  $m$  depicted in Figure 113: for  $A \in \mathbb{R}^{m \times n}$

- Given data from Example 4.2.3.1.1, for  $m=3$ ,  $n=6$ ,  $k=1$

$$A = \begin{bmatrix} -1 & 1 & 8 & 1 & 1 & 0 \\ -3 & 2 & 8 & \frac{1}{2} & \frac{1}{3} & \frac{1}{2} - \frac{1}{3} \\ -9 & 4 & 8 & \frac{1}{4} & \frac{1}{9} & \frac{1}{4} - \frac{1}{9} \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{4} \end{bmatrix} \quad (728)$$

the sparsest solution to classical linear equation  $Ax = b$  is  $x = e_4 \in \mathbb{R}^6$  (confer (741)).

Although the sparsest solution is recoverable by inspection, we discern it instead by convex iteration; namely, by iterating problem sequence (158) (530) on page 271. From the numerical data given, cardinality  $\|x\|_0 = 1$  is expected. Iteration continues until  $x^T y$  vanishes (to within some numerical precision); *id est*, until desired cardinality is achieved. But this comes not without a stall.

Stalling, whose occurrence is sensitive to initial conditions of convex iteration, is a consequence of finding a local minimum of a multimodal objective  $\langle x, y \rangle$  when regarded as simultaneously variable in  $x$  and  $y$ . (§3.14.0.0.3) Stalls are simply detected as fixed points  $x$  of infeasible cardinality, sometimes remedied by reinitializing direction vector  $y$  to a random positive state.

Bolstered by success in breaking out of a stall, we then apply convex iteration to 22,000 randomized problems:

- Given random data for  $m=3$ ,  $n=6$ ,  $k=1$ , in MATLAB notation

$$A = \text{randn}(3, 6), \quad \text{index} = \text{round}(5 * \text{rand}(1)) + 1, \quad b = \text{rand}(1) * A(:, \text{index}) \quad (848)$$

the sparsest solution  $x \propto e_{\text{index}}$  is a scaled standard basis vector.

Without convex iteration or a nonnegativity constraint  $x \succeq 0$ , rate of failure for this minimal cardinality problem  $Ax=b$  by 1-norm minimization of  $x$  is 22%. That failure rate drops to 6% with a nonnegativity constraint. If we then engage convex iteration, detect stalls, and randomly reinitialize the direction vector, failure rate drops to 0% but the amount of computation is approximately doubled.  $\square$

Stalling is not an inevitable behavior. For some problem types (beyond mere  $Ax=b$ ), convex iteration succeeds nearly all the time. Here is a cardinality problem, with noise, whose statement is just a bit more intricate but easy to solve in a few convex iterations:

#### 4.6.1.5.2 Example. *Signal dropout.*

[139, §6.2]

Signal dropout is an old problem; well studied from both an industrial and academic perspective. Essentially *dropout* means momentary loss or gap in a signal, while passing through some channel, caused by some man-made or natural phenomenon. The signal lost is assumed completely destroyed somehow. What remains within the time-gap is system or idle channel noise. The signal could be voice over Internet protocol (VoIP), for example, audio data from a *compact disc* (CD) or video data from a digital video disc (DVD), a television transmission over cable or the airwaves, or a typically ravaged cell phone communication, *etcetera*.

Here we consider signal dropout in a discrete-time signal corrupted by additive white noise assumed uncorrelated to the signal. The linear channel is assumed to introduce no filtering. We create a discretized windowed signal for this example by positively combining  $k$  randomly chosen vectors from a *discrete cosine transform* (DCT) basis denoted  $\Psi \in \mathbb{R}^{n \times n}$ . Frequency increases, in the Fourier sense, from DC toward Nyquist as column index of basis  $\Psi$  increases. Otherwise, details of the basis are unimportant except for its orthogonality  $\Psi^T = \Psi^{-1}$ . Transmitted signal is denoted

$$s = \Psi z \in \mathbb{R}^n \quad (849)$$

whose upper bound on DCT basis coefficient cardinality  $\text{card } z \leq k$  is assumed known,<sup>4.41</sup> hence a critical assumption: transmitted signal  $s$  is sparsely supported ( $k < n$ ) on the DCT basis. It is further assumed that nonzero signal coefficients in vector  $z$  place each chosen basis vector above the noise floor.

---

<sup>4.41</sup>This simplifies exposition, although it may be an unrealistic assumption in many applications.

We also assume that the gap's beginning and ending in time are precisely localized to within a sample; *id est*, index  $\ell$  locates the last sample prior to the gap's onset, while index  $n-\ell+1$  locates the first sample subsequent to the gap: for rectangularly windowed received signal  $g$  possessing a time-gap loss and additive noise  $\eta \in \mathbb{R}^n$

$$g = \begin{bmatrix} s_{1:\ell} & + & \eta_{1:\ell} \\ & & \eta_{\ell+1:n-\ell} \\ s_{n-\ell+1:n} & + & \eta_{n-\ell+1:n} \end{bmatrix} \in \mathbb{R}^n \quad (850)$$

The window is thereby centered on the gap and short enough so that the DCT spectrum of signal  $s$  can be assumed static over the window's duration  $n$ . Signal to noise ratio within this window is defined

$$\text{SNR} \triangleq 20 \log \frac{\left\| \begin{bmatrix} s_{1:\ell} \\ s_{n-\ell+1:n} \end{bmatrix} \right\|}{\|\eta\|} \quad (851)$$

In absence of noise, knowing the signal DCT basis and having a good estimate of basis coefficient cardinality makes perfectly reconstructing gap-loss easy: it amounts to solving a linear system of equations and requires little or no optimization; with *caveat*, number of equations exceeds cardinality of signal representation (roughly  $\ell \geq k$ ) with respect to DCT basis.

But addition of a significant amount of noise  $\eta$  increases level of difficulty dramatically; a 1-norm based method of reducing cardinality, for example, almost always returns DCT basis coefficients numbering in excess of minimal cardinality. We speculate that is because signal cardinality  $2\ell$  becomes the predominant cardinality. DCT basis coefficient cardinality is an explicit constraint to the optimization problem we shall pose: In presence of noise, constraints equating reconstructed signal  $f$  to received signal  $g$  are not possible. We can instead formulate the dropout recovery problem as a best approximation:

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \left\| \begin{bmatrix} f_{1:\ell} - g_{1:\ell} \\ f_{n-\ell+1:n} - g_{n-\ell+1:n} \end{bmatrix} \right\| \\ \text{subject to } & f = \Psi x \\ & x \succeq 0 \\ & \text{card } x \leq k \end{array} \quad (852)$$

We propose solving this nonconvex problem (852) by moving the cardinality constraint to the objective as a regularization term as explained in §4.6 (p.271); *id est*, by iteration of two convex problems until convergence:

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \langle x, y \rangle + \left\| \begin{bmatrix} f_{1:\ell} - g_{1:\ell} \\ f_{n-\ell+1:n} - g_{n-\ell+1:n} \end{bmatrix} \right\| \\ \text{subject to } & f = \Psi x \\ & x \succeq 0 \end{array} \quad (853)$$

and

$$\begin{array}{ll} \underset{y \in \mathbb{R}^n}{\text{minimize}} & \langle x^*, y \rangle \\ \text{subject to } & 0 \preceq y \preceq \mathbf{1} \\ & y^T \mathbf{1} = n - k \end{array} \quad (530)$$

Signal cardinality  $2\ell$  is implicit to the problem statement. When number of samples in the dropout region exceeds half the window size, then that deficient cardinality of signal remaining becomes a source of degradation to reconstruction in presence of noise. Thus, by observation, we divine a reconstruction rule for this signal dropout problem to attain good noise suppression:  $\ell$  must exceed a maximum of cardinality bounds;  $2\ell \geq \max\{2k, n/2\}$ .

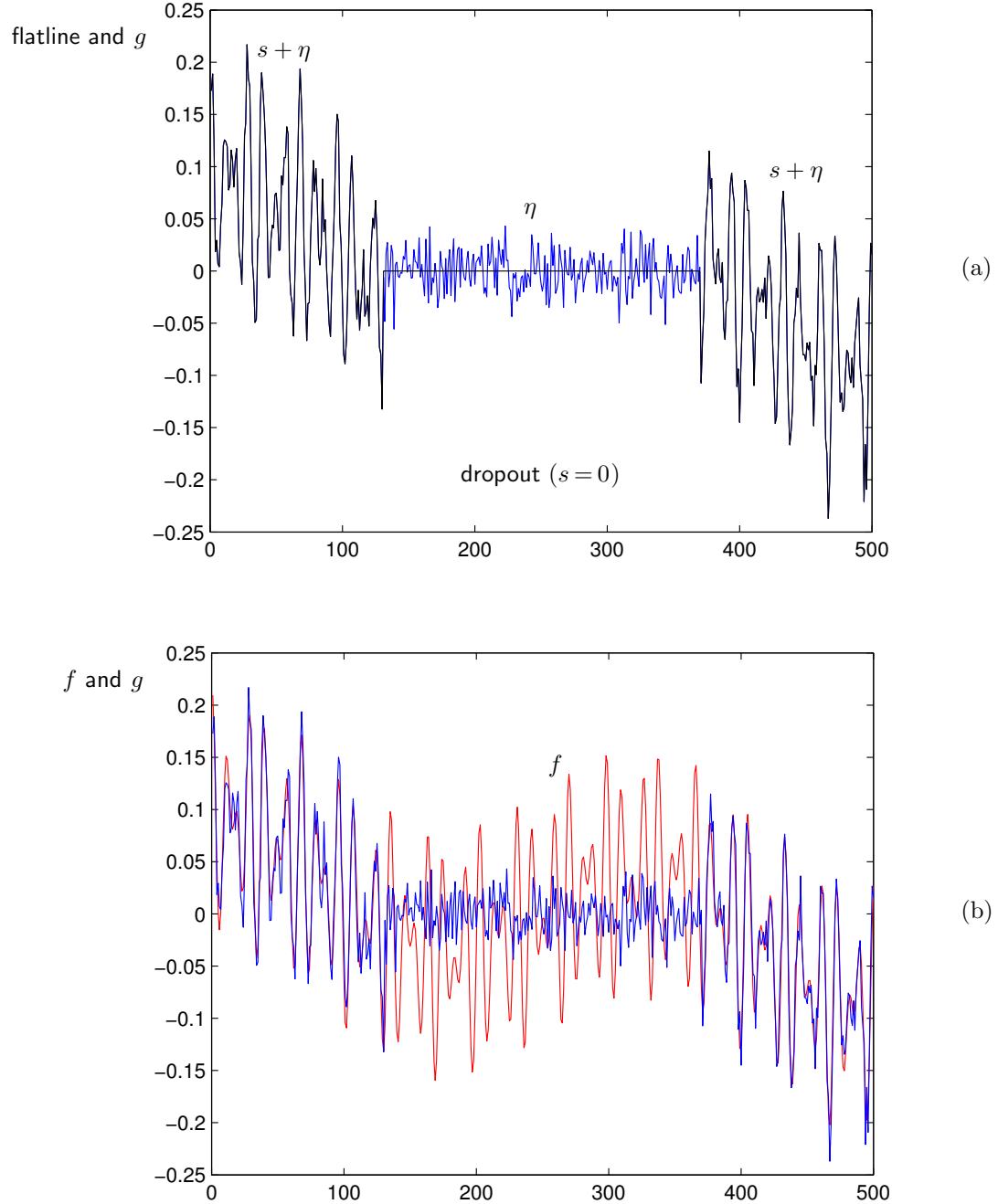


Figure 114: (a) Signal dropout in signal  $s$  corrupted by noise  $\eta$  (SNR = 10dB,  $g = s + \eta$ ). Flatline indicates duration of signal dropout. (b) Reconstructed signal  $f$  (red) overlaid with corrupted signal  $g$ .

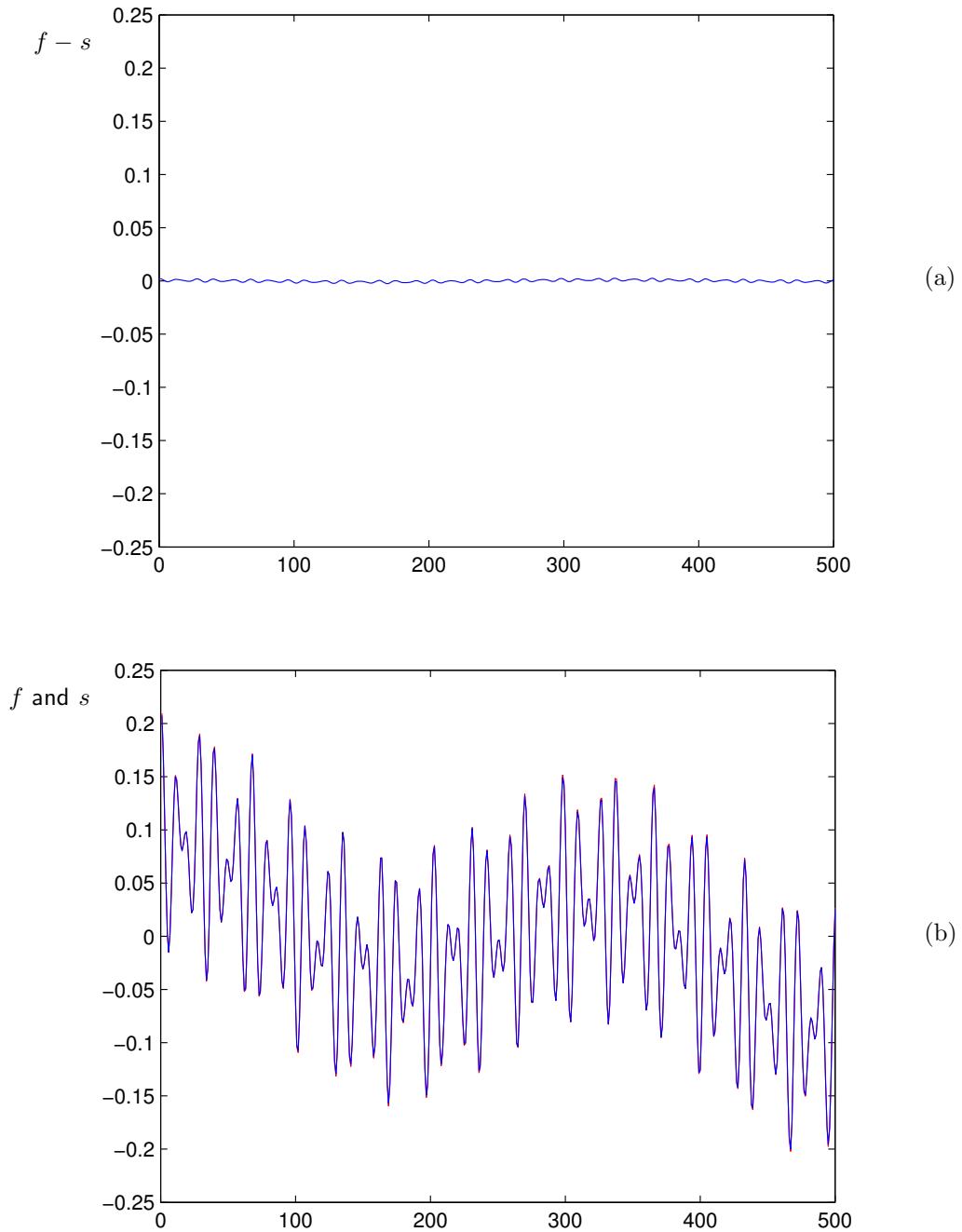


Figure 115: (a) Error signal power (reconstruction  $f$  less original noiseless signal  $s$ ) is 36dB below  $s$ . (b) Original signal  $s$  overlaid with reconstruction  $f$  (red) from signal  $g$  having dropout plus noise.

Figure 114 and Figure 115 show one realization of this dropout problem. Original signal  $s$  is created by adding four ( $k=4$ ) randomly selected DCT basis vectors, from  $\Psi$  ( $n=500$  in this example), whose amplitudes are randomly selected from a uniform distribution above the noise floor; in the interval  $[10^{-10/20}, 1]$ . Then a 240-sample dropout is realized ( $\ell=130$ ) and Gaussian noise  $\eta$  added to make corrupted signal  $g$  (from which a best approximation  $f$  will be made) having 10dB signal to noise ratio (851). The time gap contains much noise, as apparent from Figure 114a. But in only a few iterations (853) (530), original signal  $s$  is recovered with relative error power 36dB down; illustrated in Figure 115. Correct cardinality is also recovered ( $\text{card } x = \text{card } z$ ) along with the basis vector indices used to make original signal  $s$ . Approximation error is due to DCT basis coefficient estimate error. When this experiment is repeated 1000 times on noisy signals averaging 10dB SNR, the correct cardinality and indices are recovered 99% of the time with average relative error power 30dB down. Without noise, we get perfect reconstruction in one iteration. [428, MATLAB code]  $\square$

#### 4.6.1.6 Compressed sensing geometry with a nonnegative variable

It is well known that cardinality problem (535) (p.179) is easier to solve by linear programming when variable  $x$  is nonnegatively constrained than when not. We postulate a simple geometrical explanation:

Figure 74 illustrates 1-norm ball  $\mathcal{B}_1$  in  $\mathbb{R}^3$  and affine subset  $\mathcal{A}$  defined  $\{x \in \mathbb{R}^3 \mid Ax = b\}$ . Prototypical compressed sensing problem, for  $A \in \mathbb{R}^{m \times n}$

$$\begin{aligned} & \underset{x}{\text{minimize}} && \|x\|_1 \\ & \text{subject to} && Ax = b \end{aligned} \tag{523}$$

is solved when the 1-norm ball  $\mathcal{B}_1$  kisses the affine subset.

If variable  $x$  is constrained to the nonnegative orthant

$$\begin{array}{lll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \|x\|_1 & \underset{x \in \mathbb{R}^n}{\text{minimize}} & \mathbf{1}^T x & \underset{c \in \mathbb{R}, x \in \mathbb{R}^n}{\text{minimize}} & c \\ \text{subject to} & Ax = b & \equiv & Ax = b & \equiv & \text{subject to} & Ax = b \\ & x \succeq 0 & & x \succeq 0 & & & x \in c\mathcal{S} \end{array} \tag{528}$$

then 1-norm ball  $\mathcal{B}_1$  becomes nonnegative simplex  $\mathcal{S}$  in Figure 116 where

$$c\mathcal{S} = \{[I \in \mathbb{R}^{n \times n} \quad \mathbf{0} \in \mathbb{R}^n]a \mid a^T \mathbf{1} = c, a \succeq 0\} = \{x \mid x \succeq 0, \mathbf{1}^T x \leq c\} \tag{854}$$

Nonnegative simplex  $\mathcal{S}$  is the convex hull of its vertices. All  $n+1$  vertices of  $\mathcal{S}$  are constituted by standard basis vectors and the origin. In other words, all its nonzero extreme points are cardinality-1.

Affine subset  $\mathcal{A}$  kisses nonnegative simplex  $c^*\mathcal{S}$  at optimality of (528). A kissing point is achieved at  $x^*$  for optimal  $c^*$  as  $\mathcal{B}_1$  or  $\mathcal{S}$  contracts. Whereas 1-norm ball  $\mathcal{B}_1$  has only six vertices in  $\mathbb{R}^3$  corresponding to cardinality-1 solutions, simplex  $\mathcal{S}$  has three edges (along the Cartesian axes) containing an infinity of cardinality-1 solutions. And whereas  $\mathcal{B}_1$  has twelve edges containing cardinality-2 solutions,  $\mathcal{S}$  has three (out of total four) facets constituting cardinality-2 solutions. In other words, likelihood of a low-cardinality solution is higher by kissing nonnegative simplex  $\mathcal{S}$  (528) than by kissing 1-norm ball  $\mathcal{B}_1$  (523) because facial dimension (corresponding to given cardinality) is higher in  $\mathcal{S}$ .

Empirically, this observation also holds in other Euclidean dimensions; e.g., Figure 75, Figure 113.

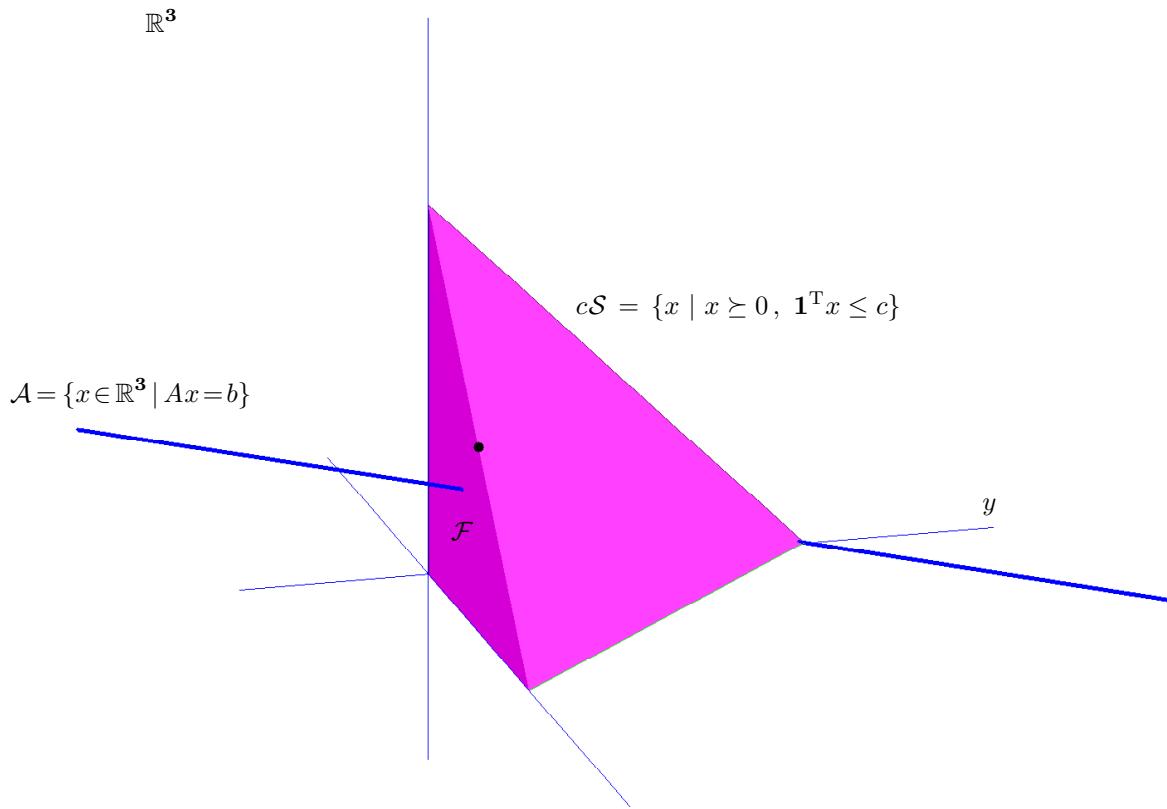


Figure 116: Simplex  $\mathcal{S}$  is convex hull of origin and all cardinality-1 nonnegative vectors of unit norm (its vertices). Line  $\mathcal{A}$ , intersecting two-dimensional (cardinality-2) face  $\mathcal{F}$  of nonnegative simplex  $c\mathcal{S}$ , emerges from  $c\mathcal{S}$  at a cardinality-1 vertex.  $\mathcal{S}$  equals nonnegative orthant  $\mathbb{R}_+^3 \cap 1\text{-norm ball } \mathcal{B}_1$  (Figure 74). Kissing point achieved when • (on edge) meets  $\mathcal{A}$  as simplex contracts (as scalar  $c$  diminishes) under optimization (528).

#### 4.6.1.7 cardinality-1 compressed sensing problem always solvable

In the special case of cardinality-1 feasible solution to nonnegative compressed sensing problem (528), there is a geometrical interpretation that leads to an algorithm.

Figure 116 illustrates a cardinality-1 feasible solution to problem (528) in  $\mathbb{R}^3$ ; a vertex solution. But *first-octant*  $\mathcal{S}$  of 1-norm ball  $\mathcal{B}_1$  does not kiss line  $\mathcal{A}$ ; which would be an optimality condition. How can we perform optimization and make  $\mathcal{A}$  intersect  $\mathcal{S}$  at a vertex? Assuming that nonnegative cardinality-1 solutions exist in the feasible set, it so happens:

##### 4.6.1.7.1 Algorithm. *Deprecation.*

Columns of *measurement matrix*  $A$ , corresponding to high cardinality solution of (528)<sup>4.42</sup> found by Simplex method [103], may be *deprecated* and the problem solved again with those columns missing. Such columns are recursively removed from  $A$  until a cardinality-1 solution is found. ¶

This algorithm intimates that either a solution to problem (528) is cardinality-1 or column indices of  $A$ , corresponding to a higher cardinality solution, do not intersect that index corresponding to a cardinality-1 feasible solution.

When problem (528) is first solved, in the example of Figure 116, solution is cardinality-2 at a kissing point on that edge of simplex  $c\mathcal{S}$  indicated by •. Imagining that the corresponding cardinality-2 face  $\mathcal{F}$  has collapsed, as a result of zeroing those two extreme points whose convex hull constructs that same edge • of  $\mathcal{F}$ , then the simplex collapses to a line segment along the  $y$  axis. When that line segment kisses  $\mathcal{A}$ , then the cardinality-1 vertex solution illustrated has been found.<sup>4.43</sup>

##### 4.6.1.7.2 Proof (pending). *Deprecation algorithm 4.6.1.7.1.*

We require proof that a cardinality-1 feasible solution to (528) cannot exist within a higher cardinality optimal solution found by Simplex method; for only then can corresponding columns of  $A$  be eliminated without precluding cardinality-1 at optimality of the deprecated problem. Crucial is the Simplex method of solution because then an optimal solution is guaranteed to reside at a vertex of the feasible set. [103, p.158] [16, p.2] ■

Although it is more efficient (compared with our algorithm) to search over individual columns of matrix  $A$  for a cardinality-1 solution known *a priori* to exist, tables are turned when cardinality exceeds 1 :

## 4.6.2 cardinality- $k$ geometric presolver

This idea of deprecating columns has foundation in convex cone theory. (§2.13.5) Removing columns (and rows)<sup>4.44</sup> from  $A \in \mathbb{R}^{m \times n}$ , in a linear program like (528) in §3.2, is known in the industry as *presolving*;<sup>4.45</sup> the elimination of redundant constraints and identically

<sup>4.42</sup>Because signed compressed sensing problem (523) can be equivalently expressed in a nonnegative variable, as we learned in Example 3.2.0.0.1 (p.177), and because a cardinality-1 constraint in (523) transforms to a cardinality-1 constraint in its nonnegative equivalent (527), then this cardinality-1 recursive reconstruction algorithm continues to hold for a signed variable as in (523).

<sup>4.43</sup>A similar argument holds for any orientation of line  $\mathcal{A}$  and cardinality-1 point of emergence from simplex  $c\mathcal{S}$ . This cardinality-1 reconstruction algorithm also holds more generally when affine subset  $\mathcal{A}$  has any higher dimension  $n - m$ .

<sup>4.44</sup>Rows of matrix  $A$  are removed based upon linear dependence. Assuming  $b \in \mathcal{R}(A)$ , corresponding entries of vector  $b$  may also be removed without loss of generality.

<sup>4.45</sup>... *presolving* can in particular do the following:

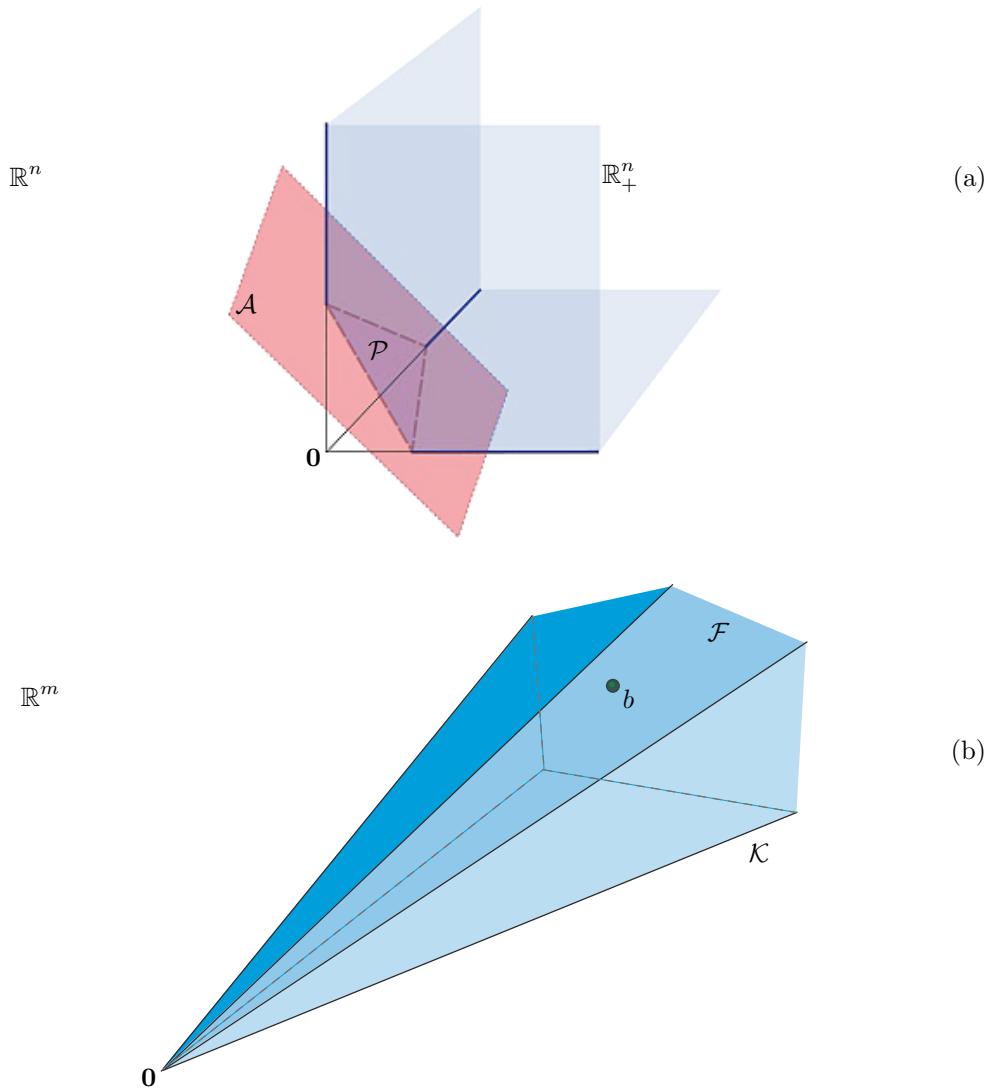


Figure 117: Constraint interpretations: (a) Halfspace-description of feasible set in problem (528) is a polyhedron  $\mathcal{P}$  formed by intersection of nonnegative orthant  $\mathbb{R}_+$  with hyperplanes  $\mathcal{A}$  prescribed by equality constraint. (Drawing by Pedro Sánchez.) (b) Vertex-description of constraints in problem (528): point  $b$  belongs to polyhedral cone  $\mathcal{K} = \{Ax \mid x \succeq 0\}$ . Number of extreme directions in  $\mathcal{K}$  may exceed dimensionality of ambient space.

zero variables prior to numerical solution. We offer a different and geometric presolver first introduced in §2.13.5:<sup>4.46</sup>

Two interpretations of the constraints from problem (528) are realized in Figure 117. Assuming that a cardinality- $k$  solution exists and matrix  $A$  describes a pointed polyhedral cone  $\mathcal{K} = \{Ax \mid x \succeq 0\}$ , as in Figure 117b, columns are removed from  $A$  if they do not belong to the smallest face  $\mathcal{F}$  of  $\mathcal{K}$  containing vector  $b$ ; those columns correspond to 0-entries in variable vector  $x$  (and *vice versa*). Generators of that smallest face always hold a minimal cardinality solution, in other words, because a generator outside the smallest face (having positive coefficient) would violate the assumption that  $b$  belongs to that face.

Benefit accrues when vector  $b$  does not belong to relative interior of  $\mathcal{K}$ ; there would be no columns to remove were  $b \in \text{rel intr } \mathcal{K}$  since the smallest face becomes cone  $\mathcal{K}$  itself (Example 4.6.2.0.2). Were  $b$  an extreme direction, at the other end of the spectrum, then the smallest face is an edge that is a ray containing  $b$ ; this geometrically describes a cardinality-1 case where all columns, save one, would be removed from  $A$ .

When vector  $b$  resides in a face  $\mathcal{F}$  of  $\mathcal{K}$  that is not cone  $\mathcal{K}$  itself, benefit is realized as a reduction in computational intensity because the consequent equivalent problem has smaller dimension. Number of columns removed depends completely on geometry of a given problem; particularly, location of  $b$  within  $\mathcal{K}$ . In the example of Figure 117b, interpreted literally in  $\mathbb{R}^3$ , all but two columns of  $A$  are discarded by our presolver when  $b$  belongs to facet  $\mathcal{F}$ .

#### 4.6.2.0.1 Exercise. Minimal cardinality generators.

Prove that generators of the smallest face  $\mathcal{F}$  of  $\mathcal{K} = \{Ax \mid x \succeq 0\}$ , containing vector  $b$ , always hold a minimal cardinality solution to  $Ax = b$ . ▼

#### 4.6.2.0.2 Example. Presolving for cardinality-2 solution to $Ax = b$ .

(confer Example 4.6.1.5.1) Again taking data from Example 4.2.3.1.1 ( $A \in \mathbb{R}^{m \times n}$ , desired cardinality of  $x$  is  $k$ ), for  $m=3$ ,  $n=6$ ,  $k=2$

$$A = \begin{bmatrix} -1 & 1 & 8 & 1 & 1 & 0 \\ -3 & 2 & 8 & \frac{1}{2} & \frac{1}{3} & \frac{1}{2} - \frac{1}{3} \\ -9 & 4 & 8 & \frac{1}{4} & \frac{1}{9} & \frac{1}{4} - \frac{1}{9} \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{4} \end{bmatrix} \quad (728)$$

proper cone  $\mathcal{K} = \{Ax \mid x \succeq 0\}$  is pointed as proven by method of §2.12.2.2. A cardinality-2 solution is known to exist; sum of the last two columns of matrix  $A$ . Generators of the smallest face that contains vector  $b$ , found by the method in Example 2.13.5.0.1, comprise the entire  $A$  matrix because  $b \in \text{intr } \mathcal{K}$  (§2.13.4.2.4). So geometry of this particular problem does not permit number of generators to be reduced below  $n$  by discerning the smallest face.<sup>4.47</sup> □

There is wondrous bonus to presolving when a constraint matrix is sparse. After columns are removed by theory of convex cones (finding the smallest face), some remaining rows may become  $\mathbf{0}^T$ , identical to other rows, or nonnegative. When nonnegative

1. Fix a variable, i.e, permanently set  $y=p$ .
2. Aggregate a variable, i.e, conclude that  $y=ax+c$  for some values  $a$  and  $c$ .
3. Multi-aggregate a variable, i.e, conclude that  $y=a_1x_1+\dots+a_kx_k+c$ .

In all cases,  $y$  will be removed from the set of “active” variables and instead added to the set of “fixed” variables. — Tobias Achterberg

<sup>4.46</sup>Comparison of computational intensity to a brute force search would pit combinatorial complexity, a binomial coefficient  $\propto \binom{n}{k}$ , against polynomial complexity of this conic presolver.

<sup>4.47</sup>But a canonical set of conically independent generators of  $\mathcal{K}$  comprise only the first two and last two columns of  $A$ .

rows appear in an equality constraint to  $\mathbf{0}$ , all nonnegative variables corresponding to nonnegative entries in those rows must vanish (§A.7.1); meaning, more columns may be removed. Once rows and columns have been removed from a constraint matrix, even more rows and columns may be removed by repeating the presolver procedure.

### 4.6.3 constraining cardinality of signed variable

Now consider a feasibility problem equivalent to the classical problem from linear algebra  $Ax = b$ , but with an upper bound  $k$  on cardinality  $\|x\|_0$ : for vector  $b \in \mathcal{R}(A)$

$$\begin{aligned} & \text{find } x \in \mathbb{R}^n \\ & \text{subject to } Ax = b \\ & \quad \|x\|_0 \leq k \end{aligned} \tag{855}$$

where  $\|x\|_0 \leq k$  means vector  $x$  has at most  $k$  nonzero entries; such a vector is presumed existent in the feasible set. Convex iteration (§4.6.1) utilizes a nonnegative variable; so absolute value  $|x|$  is needed here. We propose that nonconvex problem (855) can be equivalently written as a sequence of convex problems that move the cardinality constraint to the objective:

$$\begin{aligned} & \underset{\substack{x \in \mathbb{R}^n \\ \text{subject to } Ax = b}}{\text{minimize}} \quad \langle |x|, y \rangle \\ & \equiv \underset{\substack{x \in \mathbb{R}^n, t \in \mathbb{R}^n \\ \text{subject to } Ax = b \\ -t \preceq x \preceq t}}{\text{minimize}} \quad \langle t, y + \varepsilon \mathbf{1} \rangle \end{aligned} \tag{856}$$

$$\begin{aligned} & \underset{\substack{y \in \mathbb{R}^n \\ \text{subject to } 0 \preceq y \preceq \mathbf{1} \\ y^T \mathbf{1} = n - k}}{\text{minimize}} \quad \langle t^*, y + \varepsilon \mathbf{1} \rangle \\ & \tag{530} \end{aligned}$$

where  $\varepsilon$  is a relatively small positive constant. This sequence is iterated until a direction vector  $y$  is found that makes  $|x^*|^T y^*$  vanish. The term  $\langle t, \varepsilon \mathbf{1} \rangle$  in (856) is necessary to determine absolute value  $|x^*| = t^*$  (§3.2) because vector  $y$  can have zero-valued entries. By initializing  $y$  to  $(1-\varepsilon)\mathbf{1}$ , the first iteration of problem (856) is a 1-norm problem (519); *id est*,

$$\begin{aligned} & \underset{\substack{x \in \mathbb{R}^n, t \in \mathbb{R}^n \\ \text{subject to } Ax = b \\ -t \preceq x \preceq t}}{\text{minimize}} \quad \langle t, \mathbf{1} \rangle \\ & \equiv \underset{\substack{x \in \mathbb{R}^n \\ \text{subject to } Ax = b}}{\text{minimize}} \quad \|x\|_1 \end{aligned} \tag{523}$$

Subsequent iterations of problem (856) engaging cardinality term  $\langle t, y \rangle$  can be interpreted as corrections to this 1-norm problem leading to a 0-norm solution; vector  $y$  can be interpreted as a direction of search.

#### 4.6.3.1 local optimality

As before (§4.6.1.3), convex iteration (856) (530) always converges to a locally optimal solution; a fixed point of possibly infeasible cardinality.

#### 4.6.3.2 simple variations on a signed variable

Several useful equivalents to linear programs (856) (530) are easily devised, but their geometrical interpretation is not as apparent: *e.g.*, equivalent in the limit  $\varepsilon \rightarrow 0^+$

$$\begin{aligned} & \underset{\substack{x \in \mathbb{R}^n, t \in \mathbb{R}^n \\ \text{subject to } Ax = b \\ -t \preceq x \preceq t}}{\text{minimize}} \quad \langle t, y \rangle \\ & \tag{857} \end{aligned}$$

$$\begin{aligned} & \underset{y \in \mathbb{R}^n}{\text{minimize}} \quad \langle |x^*|, y \rangle \\ & \text{subject to} \quad 0 \leq y \leq \mathbf{1} \\ & \quad y^T \mathbf{1} = n - k \end{aligned} \tag{530}$$

We get another equivalent to linear programs (856) (530), in the limit, by interpreting problem (523) as infimum to a vertex-description of the 1-norm ball (Figure 74, Example 3.2.0.0.1, confer(522)):

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \|x\|_1 \quad \equiv \quad \underset{a \in \mathbb{R}^{2n}}{\text{minimize}} \quad \langle a, y \rangle \\ & \text{subject to} \quad Ax = b \quad \text{subject to} \quad [A \quad -A]a = b \\ & \quad a \succeq 0 \\ & \underset{y \in \mathbb{R}^{2n}}{\text{minimize}} \quad \langle a^*, y \rangle \\ & \text{subject to} \quad 0 \leq y \leq \mathbf{1} \\ & \quad y^T \mathbf{1} = 2n - k \end{aligned} \tag{530}$$

where  $x^* = [I \quad -I]a^*$ ; from which it may be rightfully construed that any vector 1-norm minimization problem has equivalent expression in a nonnegative variable.

## 4.7 Cardinality and rank constraint examples

**4.7.0.0.1 Example.** *Projection on ellipsoid boundary.* [55] [168, §5.1] [277, §2]  
Consider classical linear equation  $Ax = b$  but with constraint on norm of solution  $x$ , given matrices  $C$ , wide  $A$ , and vector  $b \in \mathcal{R}(A)$

$$\begin{aligned} & \text{find} \quad x \in \mathbb{R}^N \\ & \text{subject to} \quad Ax = b \\ & \quad \|Cx\| = 1 \end{aligned} \tag{859}$$

The set  $\{x \mid \|Cx\|=1\}$  (2) describes an ellipsoid boundary (Figure 15). This is a nonconvex problem because solution is constrained to that boundary. Assign

$$G = \begin{bmatrix} Cx \\ 1 \end{bmatrix} \begin{bmatrix} x^T C^T & 1 \end{bmatrix} = \begin{bmatrix} X & Cx \\ x^T C^T & 1 \end{bmatrix} \triangleq \begin{bmatrix} Cxx^T C^T & Cx \\ x^T C^T & 1 \end{bmatrix} \in \mathbb{S}^{N+1} \tag{860}$$

Any rank-1 solution must have this form. (§B.1.0.2) Ellipsoidally constrained feasibility problem (859) is equivalent to:

$$\begin{aligned} & \underset{X \in \mathbb{S}^N}{\text{find}} \quad x \in \mathbb{R}^N \\ & \text{subject to} \quad Ax = b \\ & \quad G = \begin{bmatrix} X & Cx \\ x^T C^T & 1 \end{bmatrix} (\succeq 0) \\ & \quad \text{rank } G = 1 \\ & \quad \text{tr } X = 1 \end{aligned} \tag{861}$$

This is transformed to an equivalent convex problem by moving the rank constraint to the objective: We iterate solution of

$$\begin{aligned} & \underset{X \in \mathbb{S}^N, x \in \mathbb{R}^N}{\text{minimize}} \quad \langle G, Y \rangle \\ & \text{subject to} \quad Ax = b \\ & \quad G = \begin{bmatrix} X & Cx \\ x^T C^T & 1 \end{bmatrix} \succeq 0 \\ & \quad \text{tr } X = 1 \end{aligned} \tag{862}$$

with

$$\begin{aligned} & \underset{Y \in \mathbb{S}^{N+1}}{\text{minimize}} \quad \langle G^*, Y \rangle \\ & \text{subject to} \quad 0 \preceq Y \preceq I \\ & \quad \text{tr } Y = N \end{aligned} \tag{863}$$

until convergence. Initially  $\mathbf{0}$ , direction matrix  $Y \in \mathbb{S}^{N+1}$  regulates rank. (1872a) Singular value decomposition  $G^* = U\Sigma Q^T \in \mathbb{S}_+^{N+1}$  ([§A.6](#)) provides a new direction matrix  $Y = U(:, 2:N+1)U(:, 2:N+1)^T$  that optimally solves (863) at each iteration. An optimal solution to (859) is thereby found in a few iterations, making convex problem (862) its equivalent.

It remains possible for the iteration to stall; were a rank-1  $G$  matrix not found. In that case, the current search direction is momentarily reversed with an added randomized element:

$$Y = -U(:, 2:N+1) * (U(:, 2:N+1)' + \text{randn}(N, 1) * U(:, 1)') \tag{864}$$

in MATLAB notation. This heuristic is quite effective for problem (859) which is exceptionally easy to solve by convex iteration.

When  $b \notin \mathcal{R}(A)$  then problem (859) must be restated as a projection:

$$\begin{aligned} & \underset{x \in \mathbb{R}^N}{\text{minimize}} \quad \|Ax - b\| \\ & \text{subject to} \quad \|Cx\| = 1 \end{aligned} \tag{865}$$

This is a projection of point  $b$  on an ellipsoid boundary because any affine transformation of an ellipsoid remains an ellipsoid. Problem (862) in turn becomes

$$\begin{aligned} & \underset{X \in \mathbb{S}^N, x \in \mathbb{R}^N}{\text{minimize}} \quad \langle G, Y \rangle + \|Ax - b\| \\ & \text{subject to} \quad G = \begin{bmatrix} X & Cx \\ x^T C^T & 1 \end{bmatrix} \succeq 0 \\ & \quad \text{tr } X = 1 \end{aligned} \tag{866}$$

We iterate this with calculation (863) of direction matrix  $Y$  as before until a rank-1  $G$  matrix is found.  $\square$

#### 4.7.0.0.2 Example. Orthonormal Procrustes.

[55]

Example 4.7.0.0.1 is extensible. An orthonormal matrix  $Q \in \mathbb{R}^{n \times p}$  is characterized  $Q^T Q = I$ . Consider the particular case  $Q = [x \ y] \in \mathbb{R}^{n \times 2}$  as variable to a Procrustes problem ([§C.3](#)): given  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{m \times 2}$ :

$$\begin{aligned} & \underset{Q \in \mathbb{R}^{n \times 2}}{\text{minimize}} \quad \|AQ - B\|_F \\ & \text{subject to} \quad Q^T Q = I \end{aligned} \tag{867}$$

which is nonconvex. By vectorizing matrix  $Q$  we can make the assignment:

$$G = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} [x^T \ y^T \ 1] = \begin{bmatrix} X & Z & x \\ Z^T & Y & y \\ x^T & y^T & 1 \end{bmatrix} \triangleq \begin{bmatrix} xx^T & xy^T & x \\ yx^T & yy^T & y \\ x^T & y^T & 1 \end{bmatrix} \in \mathbb{S}^{2n+1} \tag{868}$$

Now orthonormal Procrustes problem (867) can be equivalently restated:

$$\begin{aligned} & \underset{X, Y \in \mathbb{S}, Z, x, y}{\text{minimize}} \quad \|A[x \ y] - B\|_{\text{F}} \\ \text{subject to} \quad & G = \begin{bmatrix} X & Z & x \\ Z^T & Y & y \\ x^T & y^T & 1 \end{bmatrix} (\succeq 0) \\ & \text{rank } G = 1 \\ & \text{tr } X = 1 \\ & \text{tr } Y = 1 \\ & \text{tr } Z = 0 \end{aligned} \tag{869}$$

To solve this, we form the convex problem sequence:

$$\begin{aligned} & \underset{X, Y, Z, x, y}{\text{minimize}} \quad \|A[x \ y] - B\|_{\text{F}} + \langle G, W \rangle \\ \text{subject to} \quad & G = \begin{bmatrix} X & Z & x \\ Z^T & Y & y \\ x^T & y^T & 1 \end{bmatrix} \succeq 0 \\ & \text{tr } X = 1 \\ & \text{tr } Y = 1 \\ & \text{tr } Z = 0 \end{aligned} \tag{870}$$

and

$$\begin{aligned} & \underset{W \in \mathbb{S}^{2n+1}}{\text{minimize}} \quad \langle G^*, W \rangle \\ \text{subject to} \quad & 0 \preceq W \preceq I \\ & \text{tr } W = 2n \end{aligned} \tag{871}$$

which has an optimal solution  $W$  that is known in closed form (p.533). These two problems are iterated until convergence and a rank-1  $G$  matrix is found. A good initial value for direction matrix  $W$  is  $\mathbf{0}$ . Optimal  $Q^*$  equals  $[x^* \ y^*]$ .

Numerically, this Procrustes problem is easy to solve; a solution seems always to be found in one or few iterations. This problem formulation is extensible, of course, to orthogonal (square) matrices  $Q$ .  $\square$

#### 4.7.0.0.3 Example. Combinatorial Procrustes problem.

In case  $A, B \in \mathbb{R}^n$ , when vector  $A = \Xi B$  is known to be a permutation of vector  $B$ , solution to orthogonal Procrustes problem

$$\begin{aligned} & \underset{X \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \|A - XB\|_{\text{F}} \\ \text{subject to} \quad & X^T = X^{-1} \end{aligned} \tag{1884}$$

is not necessarily a permutation matrix  $\Xi$  even though an optimal objective value of 0 is found by the known analytical solution (§C.3). The simplest method of solution finds permutation matrix  $X^* = \Xi$  simply by sorting vector  $B$  with respect to  $A$ .

Instead of sorting, we design two different convex problems each of whose optimal solution is a permutation matrix: one design is based on rank constraint, the other on cardinality. Because permutation matrices are sparse by definition, we depart from a traditional Procrustes problem by instead demanding a vector 1-norm which is known to produce solutions more sparse than Frobenius' norm.

There are two principal facts exploited by the first convex iteration design (§4.5.1) we propose. Permutation matrices  $\Xi$  constitute:

- 1) the set of all nonnegative orthogonal matrices,
- 2) all points extreme to the polyhedron (102) of doubly stochastic matrices.

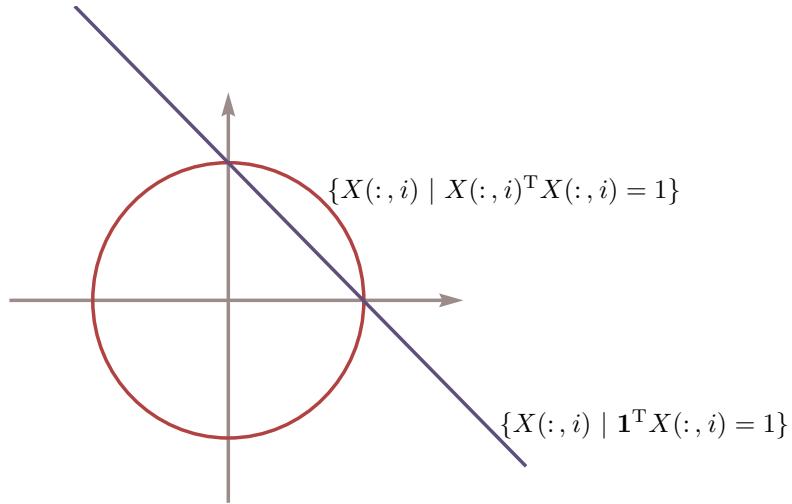


Figure 118: Permutation matrix  $i^{\text{th}}$  column-norm and column-sum constraint, abstract in two dimensions, when rank-1 constraint is satisfied. Optimal solutions reside at intersection of hyperplane with unit circle.

That means:

- 1) norm of each row and column is 1, [4.48](#)

$$\|\Xi(:, i)\| = 1, \quad \|\Xi(i, :)\| = 1, \quad i=1 \dots n \quad (872)$$

- 2) sum of each nonnegative row and column is 1, ([§2.3.2.0.4](#))

$$\Xi^T \mathbf{1} = \mathbf{1}, \quad \Xi \mathbf{1} = \mathbf{1}, \quad \Xi \geq \mathbf{0} \quad (873)$$

### solution via rank constraint

The idea is to individually constrain each column of variable matrix  $X$  to have unity norm. Matrix  $X$  must also belong to that polyhedron, (102) in the nonnegative orthant, implied by constraints (873); so each row-sum and column-sum of  $X$  must also be unity. It is this combination of nonnegativity, sum, and sum square constraints that extracts the permutation matrices: (Figure 118) given nonzero vectors  $A, B$

$$\begin{aligned} & \underset{X \in \mathbb{R}^{n \times n}, G_i \in \mathbb{S}^{n+1}}{\text{minimize}} && \|A - XB\|_1 + w \sum_{i=1}^n \langle G_i, W_i \rangle \\ & \text{subject to} && G_i = \left[ \begin{array}{cc} G_i(1:n, 1:n) & X(:, i) \\ X(:, i)^T & 1 \end{array} \right] \succeq 0 \quad \left. \right\}, \quad i=1 \dots n \end{aligned} \quad (874)$$

$$\begin{aligned} X^T \mathbf{1} &= \mathbf{1} \\ X \mathbf{1} &= \mathbf{1} \\ X &\geq \mathbf{0} \end{aligned}$$

---

[4.48](#)This fact would be superfluous were the objective of minimization linear, because the permutation matrices reside at the extreme points of a polyhedron (102) implied by (873). But as posed, only either rows or columns need be constrained to unit norm because matrix orthogonality implies transpose orthogonality. ([§B.5.2](#)) Absence of vanishing inner product constraints that help define orthogonality, like  $\text{tr } Z = 0$  from Example 4.7.0.0.2, is a consequence of nonnegativity; *id est*, the only orthogonal matrices having exclusively nonnegative entries are permutations of the Identity.

where  $w \approx 10$  positively weights the rank regularization term. Optimal solutions  $G_i^*$  are key to finding direction matrices  $W_i$  for the next iteration of semidefinite programs (874) (875):

$$\begin{array}{ll} \text{minimize}_{\substack{W_i \in \mathbb{S}^{n+1}}} & \langle G_i^*, W_i \rangle \\ \text{subject to} & \left. \begin{array}{l} 0 \preceq W_i \preceq I \\ \text{tr } W_i = n \end{array} \right\}, \quad i=1 \dots n \end{array} \quad (875)$$

Direction matrices thus found lead toward rank-1 matrices  $G_i^*$  on subsequent iterations. Constraint on trace of  $G_i^*$  normalizes the  $i^{\text{th}}$  column of  $X^*$  to unity because (*confer* p.359)

$$G_i^* = \begin{bmatrix} X^*(:, i) \\ 1 \end{bmatrix} [X^*(:, i)^T \ 1] \quad (876)$$

at convergence. Binary-valued  $X^*$  column entries result from the further sum constraint  $X\mathbf{1}=\mathbf{1}$ . Columnar orthogonality is a consequence of the further transpose-sum constraint  $X^T\mathbf{1}=\mathbf{1}$  in conjunction with nonnegativity constraint  $X \geq \mathbf{0}$ ; but we leave proof of orthogonality an exercise. The optimal objective value is 0 for both semidefinite programs when vectors  $A$  and  $B$  are related by permutation. In any case, optimal solution  $X^*$  becomes a permutation matrix  $\Xi$ .

Because there are  $n$  direction matrices  $W_i$  to find, it can be advantageous to invoke a known closed-form solution for each from page 533. What makes this combinatorial problem more tractable are relatively small semidefinite constraints in (874). (*confer*(870)) When a permutation  $A$  of vector  $B$  exists, number of iterations can be as small as 1. But this combinatorial Procrustes problem can be made even more challenging when vector  $A$  has repeated entries.

#### solution via cardinality constraint

Now the idea is to force solution at a vertex of permutation polyhedron (102) by finding a solution of desired sparsity. Because permutation matrix  $X$  is  $n$ -sparse by assumption, this combinatorial Procrustes problem may instead be formulated as a compressed sensing problem with convex iteration on cardinality of vectorized  $X$  (§4.6.1): given nonzero vectors  $A, B$

$$\begin{array}{ll} \text{minimize}_{\substack{X \in \mathbb{R}^{n \times n}}} & \|A - XB\|_1 + w\langle X, Y \rangle \\ \text{subject to} & \begin{array}{l} X^T\mathbf{1} = \mathbf{1} \\ X\mathbf{1} = \mathbf{1} \\ X \geq \mathbf{0} \end{array} \end{array} \quad (877)$$

where direction vector  $Y$  is an optimal solution to

$$\begin{array}{ll} \text{minimize}_{\substack{Y \in \mathbb{R}^{n \times n}}} & \langle X^*, Y \rangle \\ \text{subject to} & \begin{array}{l} \mathbf{0} \leq Y \leq \mathbf{1} \\ \mathbf{1}^T Y \mathbf{1} = n^2 - n \end{array} \end{array} \quad (530)$$

each a linear program. In this circumstance, use of closed-form solution for direction vector  $Y$  is discouraged. When vector  $A$  is a permutation of  $B$ , both linear programs have objectives that converge to 0. When vectors  $A$  and  $B$  are permutations and no entries of  $A$  are repeated, optimal solution  $X^*$  can be found as soon as the first iteration.

In any case,  $X^* = \Xi$  is a permutation matrix. □

#### 4.7.0.0.4 Exercise. Combinatorial Procrustes constraints.

Assume that the objective of semidefinite program (874) is 0 at optimality. Prove that the constraints in program (874) are necessary and sufficient to produce a permutation matrix as optimal solution. Alternatively and equivalently, prove those constraints necessary and sufficient to optimally produce a nonnegative orthogonal matrix. ▼

**4.7.0.0.5 Example.** *Tractable polynomial constraint.*

The set of all coefficients for which a multivariate polynomial were convex is generally difficult to determine. But the ability to handle rank constraints makes any nonconvex polynomial constraint transformable to a convex constraint. *All optimization problems having polynomial objective and polynomial constraints can be reformulated as a semidefinite program with a rank-1 constraint.* [317] Suppose we require

$$3 + 2x - xy \leq 0 \quad (878)$$

Identify

$$G = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \begin{bmatrix} x & y & 1 \end{bmatrix} = \begin{bmatrix} x^2 & xy & x \\ xy & y^2 & y \\ x & y & 1 \end{bmatrix} \in \mathbb{S}^3 \quad (879)$$

Then nonconvex polynomial constraint (878) is equivalent to constraint set

$$\begin{aligned} \text{tr}(GA) &\leq 0 \\ G_{33} &= 1 \\ (G \succeq 0) \\ \text{rank } G &= 1 \end{aligned} \quad (880)$$

with direct correspondence to sense of trace inequality where  $G$  is assumed symmetric (§B.1.0.2) and

$$A = \begin{bmatrix} 0 & -\frac{1}{2} & 1 \\ -\frac{1}{2} & 0 & 0 \\ 1 & 0 & 3 \end{bmatrix} \in \mathbb{S}^3 \quad (881)$$

Then the method of convex iteration from §4.5.1 is applied to implement the rank constraint.  $\square$

**4.7.0.0.6 Exercise.** *Binary Pythagorean theorem.*

The technique in Example 4.7.0.0.5 is extensible to any quadratic constraint; e.g.,  $x^T A x + 2b^T x + c \leq 0$ ,  $x^T A x + 2b^T x + c \geq 0$ , and  $x^T A x + 2b^T x + c = 0$ . Write a rank-constrained semidefinite program to solve (Figure 118)

$$\left\{ \begin{array}{l} x + y = 1 \\ x^2 + y^2 = 1 \end{array} \right. \quad (882)$$

whose feasible set is not connected. **cvx** is a high-level prototyping language [191] for Optimization that runs under MATLAB. Implement this system in **cvx** by convex iteration.

**4.7.0.0.7 Example.** *High order polynomials.*

Consider nonconvex problem from Canadian Mathematical Olympiad 1999:

$$\begin{array}{ll} \text{find} & x, y, z \\ \text{subject to} & x^2 y + y^2 z + z^2 x = \frac{2^2}{3^3} \\ & x + y + z = 1 \\ & x, y, z \geq 0 \end{array} \quad (883)$$

We wish to solve for, what is known to be, a tight upper bound  $\frac{2^2}{3^3}$  on the constrained polynomial  $x^2 y + y^2 z + z^2 x$  by transformation to a rank-constrained semidefinite program.

First identify

$$G = \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} [x \ y \ z \ 1] = \begin{bmatrix} x^2 & xy & zx & x \\ xy & y^2 & yz & y \\ zx & yz & z^2 & z \\ x & y & z & 1 \end{bmatrix} \in \mathbb{S}^4 \quad (884)$$

$$X = \begin{bmatrix} x^2 \\ y^2 \\ z^2 \\ x \\ y \\ z \\ 1 \end{bmatrix} [x^2 \ y^2 \ z^2 \ x \ y \ z \ 1] = \begin{bmatrix} x^4 & x^2y^2 & z^2x^2 & x^3 & x^2y & zx^2 & x^2 \\ x^2y^2 & y^4 & y^2z^2 & xy^2 & y^3 & y^2z & y^2 \\ z^2x^2 & y^2z^2 & z^4 & z^2x & yz^2 & z^3 & z^2 \\ x^3 & xy^2 & z^2x & x^2 & xy & zx & x \\ x^2y & y^3 & yz^2 & xy & y^2 & yz & y \\ zx^2 & y^2z & z^3 & zx & yz & z^2 & z \\ x^2 & y^2 & z^2 & x & y & z & 1 \end{bmatrix} \in \mathbb{S}^7 \quad (885)$$

then apply convex iteration (§4.5.1) to implement rank constraints:

$$\begin{array}{ll} \text{find} & b \\ \text{subject to} & \text{tr}(XE) = \frac{2^2}{3^3} \\ & G = \begin{bmatrix} A & b \\ b^T & 1 \end{bmatrix} (\succeq 0) \\ & X = \begin{bmatrix} C & \begin{bmatrix} \delta(A) \\ b \end{bmatrix} \\ \begin{bmatrix} \delta(A)^T & b^T \end{bmatrix} & 1 \end{bmatrix} (\succeq 0) \\ & \mathbf{1}^T b = 1 \\ & b \succeq 0 \\ & \text{rank } G = 1 \\ & \text{rank } X = 1 \end{array} \quad (886)$$

where

$$E = \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \frac{1}{2} \in \mathbb{S}^7 \quad (887)$$

[422, MATLAB code]. Positive semidefiniteness is optional only when rank-1 constraints are explicit by Theorem A.3.1.0.7. Optimal solution  $(x, y, z) = (0, \frac{2}{3}, \frac{1}{3})$  to problem (883) is not unique.  $\square$

#### 4.7.0.0.8 Exercise. Motzkin polynomial.

Prove  $xy^2 + x^2y - 3xy + 1$  to be nonnegative on the nonnegative orthant.  $\blacktriangledown$

#### 4.7.0.0.9 Example. Boolean vector satisfying $Ax \leq b$ . (confer §4.2.3.1.1)

Now we consider solution to a discrete problem whose only known analytical method of solution is combinatorial in complexity: given  $A \in \mathbb{R}^{M \times N}$  and  $b \in \mathbb{R}^M$

$$\begin{array}{ll} \text{find} & x \in \mathbb{R}^N \\ \text{subject to} & Ax \leq b \\ & \delta(xx^T) = \mathbf{1} \end{array} \quad (888)$$

This nonconvex problem demands a Boolean solution [ $x_i = \pm 1$ ,  $i = 1 \dots N$ ].

Assign a rank-1 matrix of variables; symmetric variable matrix  $X$  and solution vector  $x$ :

$$G = \begin{bmatrix} x \\ 1 \end{bmatrix} \begin{bmatrix} x^T & 1 \end{bmatrix} = \begin{bmatrix} X & x \\ x^T & 1 \end{bmatrix} \triangleq \begin{bmatrix} xx^T & x \\ x^T & 1 \end{bmatrix} \in \mathbb{S}^{N+1} \quad (889)$$

Then design an equivalent semidefinite feasibility problem to find a Boolean solution to  $Ax \preceq b$ :

$$\begin{array}{ll} \text{find} & x \in \mathbb{R}^N \\ \text{subject to} & Ax \preceq b \\ & G = \begin{bmatrix} X & x \\ x^T & 1 \end{bmatrix} (\succeq 0) \\ & \text{rank } G = 1 \\ & \delta(X) = \mathbf{1} \end{array} \quad (890)$$

where  $x_i^* \in \{-1, 1\}$ ,  $i = 1 \dots N$ . The two variables  $X$  and  $x$  are made dependent via their assignment to rank-1 matrix  $G$ . By (1779), an optimal rank-1 matrix  $G^*$  must take the form (889).

As before, we regularize the rank constraint by introducing a direction matrix  $Y$  into the objective:

$$\begin{array}{ll} \text{minimize} & \langle G, Y \rangle \\ \text{subject to} & Ax \preceq b \\ & G = \begin{bmatrix} X & x \\ x^T & 1 \end{bmatrix} \succeq 0 \\ & \delta(X) = \mathbf{1} \end{array} \quad (891)$$

Solution of this semidefinite program is iterated with calculation of the direction matrix  $Y$  from semidefinite program (863). At convergence, in the sense (802), convex problem (891) becomes equivalent to nonconvex Boolean problem (888).

Direction matrix  $Y$  can be an orthogonal projector having closed-form expression, by (1872a), although convex iteration is not a projection method. (§4.5.1.1) Given randomized data  $A$  and  $b$  for a large problem, we find that stalling becomes likely (convergence of the iteration to a positive objective  $\langle G^*, Y \rangle$ ). To overcome this behavior, we introduce a heuristic into the implementation on *Wikimization* [411] that momentarily reverses direction of search (like (864)) upon stall detection. We find that rate of convergence can be sped significantly by detecting stalls early.  $\square$

#### 4.7.0.0.10 Example. Variable-vector normalization.

Suppose, within some convex optimization problem, we want vector variables  $x, y \in \mathbb{R}^N$  constrained by a nonconvex equality:

$$x \|y\| = y \quad (892)$$

*id est*,  $\|x\| = 1$  and  $x$  points in the same direction as  $y \neq \mathbf{0}$ ; e.g,

$$\begin{array}{ll} \text{minimize} & f(x, y) \\ \text{subject to} & (x, y) \in \mathcal{C} \\ & x \|y\| = y \end{array} \quad (893)$$

where  $f$  is some convex function and  $\mathcal{C}$  is some convex set. We can realize the nonconvex equality by constraining rank and adding a regularization term to the objective. Make the

assignment:

$$G = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \begin{bmatrix} x^T & y^T & 1 \end{bmatrix} = \begin{bmatrix} X & Z & x \\ Z & Y & y \\ x^T & y^T & 1 \end{bmatrix} \triangleq \begin{bmatrix} xx^T & xy^T & x \\ yx^T & yy^T & y \\ x^T & y^T & 1 \end{bmatrix} \in \mathbb{S}^{2N+1} \quad (894)$$

where  $X, Y \in \mathbb{S}^N$ , also  $Z \in \mathbb{S}^N$  [sic]. Any rank-1 solution must take the form of (894). (§B.1) The problem statement equivalent to (893) is then written

$$\begin{aligned} & \underset{X, Y \in \mathbb{S}, Z, x, y}{\text{minimize}} && f(x, y) + \|X - Y\|_F \\ & \text{subject to} && (x, y) \in \mathcal{C} \\ & && G = \begin{bmatrix} X & Z & x \\ Z & Y & y \\ x^T & y^T & 1 \end{bmatrix} (\succeq 0) \\ & && \text{rank } G = 1 \\ & && \text{tr}(X) = 1 \\ & && \delta(Z) \succeq 0 \end{aligned} \quad (895)$$

The trace constraint on  $X$  normalizes vector  $x$  while the diagonal constraint on  $Z$  maintains sign between respective entries of  $x$  and  $y$ . Regularization term  $\|X - Y\|_F$  then makes  $x$  equal to  $y$  to within a real scalar; (§C.2.0.0.2) in this case, a positive scalar. To make this program solvable by convex iteration, as explained in Example 4.5.1.2.4 and other previous examples, we move the rank constraint to the objective

$$\begin{aligned} & \underset{X, Y, Z, x, y}{\text{minimize}} && f(x, y) + \|X - Y\|_F + \langle G, W \rangle \\ & \text{subject to} && (x, y) \in \mathcal{C} \\ & && G = \begin{bmatrix} X & Z & x \\ Z & Y & y \\ x^T & y^T & 1 \end{bmatrix} \succeq 0 \\ & && \text{tr}(X) = 1 \\ & && \delta(Z) \succeq 0 \end{aligned} \quad (896)$$

by introducing a direction matrix  $W$  found from (1872a):

$$\begin{aligned} & \underset{W \in \mathbb{S}^{2N+1}}{\text{minimize}} && \langle G^*, W \rangle \\ & \text{subject to} && 0 \preceq W \preceq I \\ & && \text{tr } W = 2N \end{aligned} \quad (897)$$

This semidefinite program has an optimal solution that is known in closed form. Iteration (896) (897) terminates when  $\text{rank } G = 1$  and linear regularization  $\langle G, W \rangle$  vanishes to within some numerical tolerance in (896); typically, in two iterations. If function  $f$  competes too much with the regularization, positively weighting each regularization term will become required. At convergence, problem (896) becomes a convex equivalent to the original nonconvex problem (893).  $\square$

#### 4.7.0.0.11 Example. FAST MAX CUT.

[126]

*Let  $\Gamma$  be an  $n$ -node graph, and let the arcs  $(i, j)$  of the graph be associated with ... weights  $a_{ij}$ . The problem is to find a cut of the largest possible weight, i.e., to partition the set of nodes into two parts  $\mathcal{M}_c, \mathcal{M}'_c$  in such a way that the total weight of all arcs linking  $\mathcal{M}_c$  and  $\mathcal{M}'_c$  (i.e., with one incident node in  $\mathcal{M}_c$  and the other one in  $\mathcal{M}'_c$  [Figure 119]) is as large as possible.*

—[35, §4.3.3]

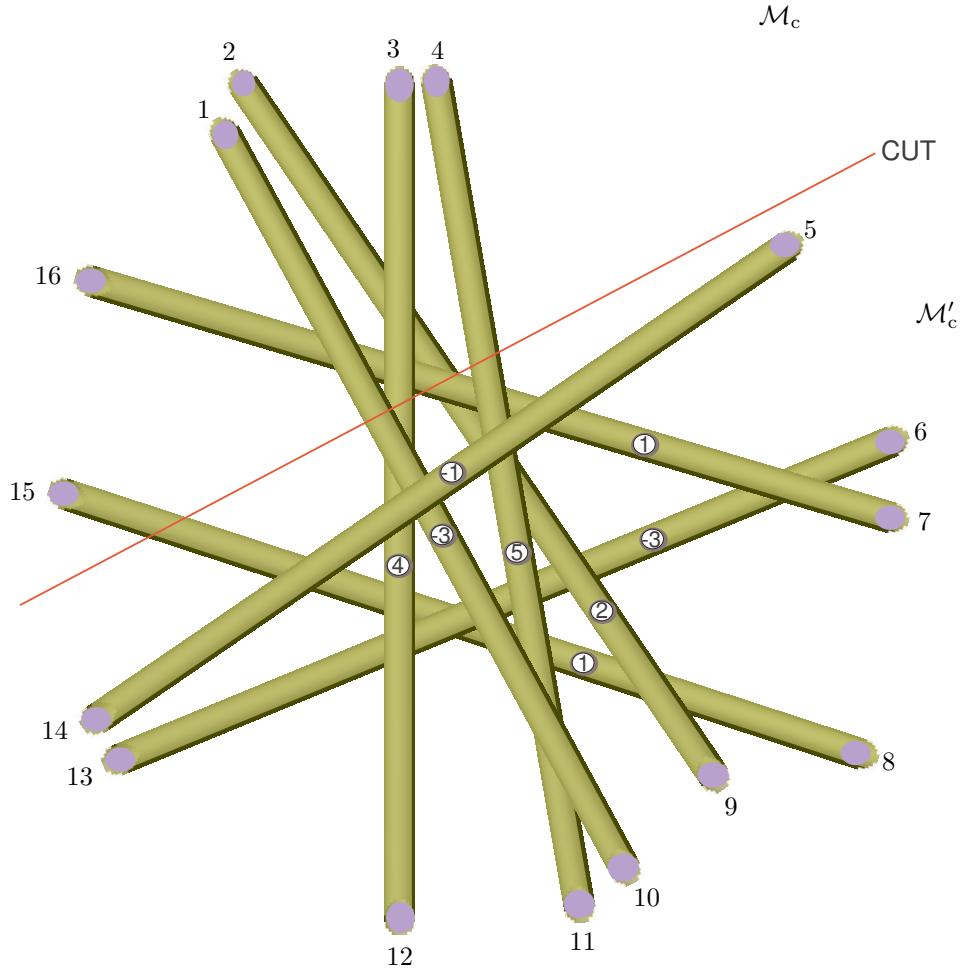


Figure 119: A CUT partitions nodes  $\{i=1 \dots 16\}$  of this graph into  $\mathcal{M}_c$  and  $\mathcal{M}'_c$ . Linear arcs have circled weights. The problem is to find a cut maximizing total weight of all arcs linking partitions made by the cut.

Literature on the MAX CUT problem is vast because this problem has elegant primal and dual formulation, its solution is very difficult, and there exist many commercial applications; *e.g.*, semiconductor design [143], quantum computing [449].

Our purpose here is to demonstrate how iteration of two simple convex problems can quickly converge to an optimal solution of the MAX CUT problem with a 98% success rate, on average.<sup>4.49</sup> MAX CUT is stated:

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{maximize}} \quad \sum_{1 \leq i < j \leq n} a_{ij} (1 - x_i x_j)^{\frac{1}{2}} \\ & \text{subject to} \quad \delta(x x^T) = \mathbf{1} \end{aligned} \tag{898}$$

where  $[a_{ij}]$  are real arc weights, and vector  $x = [x_i] \in \mathbb{R}^n$  corresponds to the  $n$  nodes; specifically,

$$\begin{aligned} \text{node } i \in \mathcal{M}_c & \Leftrightarrow x_i = 1 \\ \text{node } i \in \mathcal{M}'_c & \Leftrightarrow x_i = -1 \end{aligned} \tag{899}$$

<sup>4.49</sup>We term our solution to MAX CUT *fast* because we sacrifice a little accuracy to achieve speed; *id est*, only about two or three convex iterations, achieved by heavily weighting a rank regularization term.

If nodes  $i$  and  $j$  have the same binary value  $x_i$  and  $x_j$ , then they belong to the same partition and contribute nothing to the cut. Arc  $(i, j)$  traverses the cut, otherwise, adding its weight  $a_{ij}$  to the cut.

MAX CUT statement (898) is the same as, for  $A = [a_{ij}] \in \mathbb{S}^n$

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{maximize}} \quad \frac{1}{4} \langle \mathbf{1}\mathbf{1}^T - xx^T, A \rangle \\ & \text{subject to} \quad \delta(xx^T) = \mathbf{1} \end{aligned} \quad (900)$$

Because of Boolean assumption  $\delta(xx^T) = \mathbf{1}$

$$\langle \mathbf{1}\mathbf{1}^T - xx^T, A \rangle = \langle xx^T, \delta(A\mathbf{1}) - A \rangle \quad (901)$$

so problem (900) is the same as

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{maximize}} \quad \frac{1}{4} \langle xx^T, \delta(A\mathbf{1}) - A \rangle \\ & \text{subject to} \quad \delta(xx^T) = \mathbf{1} \end{aligned} \quad (902)$$

This MAX CUT problem is combinatorial (nonconvex).

Because an estimate of upper bound to MAX CUT is needed to ascertain convergence when vector  $x$  has large dimension, we digress to derive the dual problem: Directly from (902), its Lagrangian is [65, §5.1.5] (1572)

$$\begin{aligned} \mathfrak{L}(x, \nu) &= \frac{1}{4} \langle xx^T, \delta(A\mathbf{1}) - A \rangle + \langle \nu, \delta(xx^T) - \mathbf{1} \rangle \\ &= \frac{1}{4} \langle xx^T, \delta(A\mathbf{1}) - A \rangle + \langle \delta(\nu), xx^T \rangle - \langle \nu, \mathbf{1} \rangle \\ &= \frac{1}{4} \langle xx^T, \delta(A\mathbf{1} + 4\nu) - A \rangle - \langle \nu, \mathbf{1} \rangle \end{aligned} \quad (903)$$

where quadratic  $x^T(\delta(A\mathbf{1} + 4\nu) - A)x$  has supremum 0 if  $\delta(A\mathbf{1} + 4\nu) - A$  is assumed negative semidefinite, and has supremum  $\infty$  otherwise. The finite supremum

$$g(\nu) = \sup_{x \in \mathbb{R}^n} \mathfrak{L}(x, \nu) = \begin{cases} -\langle \nu, \mathbf{1} \rangle, & \text{assuming } A - \delta(A\mathbf{1} + 4\nu) \succeq 0 \\ \infty & \text{otherwise} \end{cases} \quad (904)$$

is chosen as the objective of minimization to dual (convex semidefinite) problem

$$\begin{aligned} & \underset{\nu \in \mathbb{R}^n}{\text{minimize}} \quad -\nu^T \mathbf{1} \\ & \text{subject to} \quad A - \delta(A\mathbf{1} + 4\nu) \succeq 0 \end{aligned} \quad (905)$$

whose optimal value  $(-\nu^T \mathbf{1})$  provides an upper bound to MAX CUT but is not tight<sup>4.50</sup> ( $\frac{1}{4} \langle xx^T, \delta(A\mathbf{1}) - A \rangle < g(\nu)$ , duality gap is nonzero); [178] problem (905) is not a strong dual to (902).<sup>4.51</sup>

To transform MAX CUT to its convex equivalent, first define

$$X = xx^T \in \mathbb{S}^n \quad (910)$$

then MAX CUT (902) becomes

$$\begin{aligned} & \underset{X \in \mathbb{S}^n}{\text{maximize}} \quad \frac{1}{4} \langle X, \delta(A\mathbf{1}) - A \rangle \\ & \text{subject to} \quad \delta(X) = \mathbf{1} \\ & \quad (X \succeq 0) \\ & \quad \text{rank } X = 1 \end{aligned} \quad (906)$$

<sup>4.50</sup>Taking the dual of dual problem (905) would provide (906) but without the rank constraint. [171] Dual of a dual of even a convex primal problem is not necessarily the same primal problem; although, optimal solution of one can be obtained from the other.

<sup>4.51</sup>Even so, empirically, binary solution  $\arg \sup_{x \in \mathbb{B}_{\pm}^n} \mathfrak{L}(x, \nu^*)$  to (903) is optimal to (902).

whose rank constraint can be regularized as in

$$\begin{aligned} & \underset{X \in \mathbb{S}^n}{\text{maximize}} \quad \frac{1}{4} \langle X, \delta(A\mathbf{1}) - A \rangle - w \langle X, W \rangle \\ & \text{subject to} \quad \delta(X) = \mathbf{1} \\ & \quad X \succeq 0 \end{aligned} \tag{907}$$

where  $w \approx 1000$  is a nonnegative fixed weight, and  $W$  is a direction matrix determined from

$$\begin{aligned} \sum_{i=2}^n \lambda(X^*)_i &= \underset{W \in \mathbb{S}^n}{\text{minimize}} \quad \langle X^*, W \rangle \\ &\text{subject to} \quad 0 \preceq W \preceq I \\ & \quad \text{tr } W = n - 1 \end{aligned} \tag{1872a}$$

which has an optimal solution that is known in closed form. These two problems (907) and (1872a) are iterated until convergence as defined on page 248.

Because convex problem statement (907) is so elegant, it is numerically solvable for large binary vectors within reasonable time.<sup>4.52</sup> To test our convex iterative method, we compare an optimal convex result to an actual solution of the MAX CUT problem found by performing a brute force combinatorial search of (902)<sup>4.53</sup> for a tight upper bound. Search-time limits binary vector lengths to 24 bits (about five days CPU time). 98% accuracy, actually obtained, is independent of binary vector length (12, 13, 20, 24) when averaged over more than 231 problem instances including planar, randomized, and toroidal graphs.<sup>4.54</sup> When failure occurred, large and small errors were manifest. That same 98% average accuracy is presumed maintained when binary vector length is further increased. A MATLAB program is provided on *Wikimization* [417].  $\square$

#### 4.7.0.0.12 Example. Cardinality/rank problem.

d'Aspremont, El Ghaoui, Jordan, & Lanckriet [104] propose approximating a positive semidefinite matrix  $A \in \mathbb{S}_+^N$  by a rank-one matrix having constraint on cardinality  $c$ : for  $0 < c < N$

$$\begin{aligned} & \underset{z}{\text{minimize}} \quad \|A - zz^T\|_F \\ & \text{subject to} \quad \text{card } z \leq c \end{aligned} \tag{908}$$

which, they explain, is a hard problem equivalent to

$$\begin{aligned} & \underset{x}{\text{maximize}} \quad x^T A x \\ & \text{subject to} \quad \|x\| = 1 \\ & \quad \text{card } x \leq c \end{aligned} \tag{909}$$

where  $z \triangleq \sqrt{\lambda} x$  and where optimal solution  $x^*$  is a *principal eigenvector* (1865) ( $\S$ A.5) of  $A$  and  $\lambda = x^{*\top} A x^*$  is the *principal eigenvalue* [181, p.331] when  $c$  is true cardinality of that eigenvector. This is *principal component analysis* with a cardinality constraint which controls solution sparsity. Define the matrix variable

$$X \triangleq xx^T \in \mathbb{S}^N \tag{910}$$

<sup>4.52</sup>We solved for a length-250 binary vector in only a few minutes and convex iterations on a 2006 vintage laptop Core 2 CPU (Intel T7400@2.16GHz, 666MHz FSB).

<sup>4.53</sup>more computationally intensive than the proposed convex iteration by many orders of magnitude. Solving MAX CUT by searching over all binary vectors of length 100, for example, would occupy a contemporary supercomputer for a million years.

<sup>4.54</sup>Existence of a polynomial-time approximation to MAX CUT with accuracy provably better than 94.11% would refute NP-hardness; which Håstad believes to be highly unlikely. [207, thm.8.2] [208]

whose desired rank is 1, and whose desired diagonal cardinality

$$\text{card } \delta(X) \equiv \text{card } x \quad (911)$$

is equivalent to cardinality  $c$  of vector  $x$ . Then we can transform cardinality problem (909) to an equivalent in new variable  $X$ :<sup>4.55</sup>

$$\begin{aligned} & \underset{X \in \mathbb{S}^N}{\text{maximize}} \quad \langle X, A \rangle \\ & \text{subject to} \quad \langle X, I \rangle = 1 \\ & \quad (X \succeq 0) \\ & \quad \text{rank } X = 1 \\ & \quad \text{card } \delta(X) \leq c \end{aligned} \quad (912)$$

We transform problem (912) to an equivalent convex problem by introducing two direction matrices into regularization terms:  $W$  to achieve desired cardinality  $\text{card } \delta(X)$ , and  $Y$  to find an approximating rank-one matrix  $X$ :

$$\begin{aligned} & \underset{X \in \mathbb{S}^N}{\text{maximize}} \quad \langle X, A - w_1 Y \rangle - w_2 \langle \delta(X), \delta(W) \rangle \\ & \text{subject to} \quad \langle X, I \rangle = 1 \\ & \quad X \succeq 0 \end{aligned} \quad (913)$$

where  $w_1$  and  $w_2$  are positive scalars respectively weighting  $\text{tr}(XY)$  and  $\delta(X)^T \delta(W)$  just enough to insure that they vanish to within some numerical precision, where direction matrix  $Y$  is an optimal solution to semidefinite program

$$\begin{aligned} & \underset{Y \in \mathbb{S}^N}{\text{minimize}} \quad \langle X^*, Y \rangle \\ & \text{subject to} \quad 0 \preceq Y \preceq I \\ & \quad \text{tr } Y = N - 1 \end{aligned} \quad (914)$$

and where diagonal direction matrix  $W \in \mathbb{S}^N$  optimally solves linear program

$$\begin{aligned} & \underset{W=\delta^2(Y)}{\text{minimize}} \quad \langle \delta(X^*), \delta(W) \rangle \\ & \text{subject to} \quad 0 \preceq \delta(W) \preceq \mathbf{1} \\ & \quad \text{tr } W = N - c \end{aligned} \quad (915)$$

Both direction matrix programs are derived from (1872a) whose analytical solution is known but is not necessarily unique. We emphasize (*confer p.248*): because this iteration (913) (914) (915) (initial  $Y, W = \mathbf{0}$ ) is not a projection method (§4.5.1.1), success relies on existence of matrices in the feasible set of (913) having desired rank and diagonal cardinality. In particular, the feasible set of convex problem (913) is a Fantope (92) whose extreme points constitute the set of all normalized rank-one matrices; among those are found rank-one matrices of any desired diagonal cardinality.

Convex problem (913) is neither a relaxation of cardinality problem (909); instead, problem (913) becomes a convex equivalent to (909) at global optimality of iteration (913) (914) (915). Because the feasible set of problem (913) contains all normalized rank-one (§B.1) symmetric matrices of every nonzero diagonal cardinality, a constraint too low or high in cardinality  $c$  will not prevent solution. An optimal rank-one solution  $X^*$ , whose diagonal cardinality is equal to cardinality of a principal eigenvector of matrix  $A$ , will produce the least residual Frobenius norm (to within machine noise processes) in the original problem statement (908).  $\square$

---

<sup>4.55</sup>A semidefiniteness constraint  $X \succeq 0$  is not required, theoretically, because positive semidefiniteness of a rank-1 matrix is enforced by symmetry. (Theorem A.3.1.0.7)



Figure 120: Shepp-Logan phantom from MATLAB *image processing toolbox*.

#### 4.7.0.0.13 Example. Compressive sampling of a phantom.

In Summer 2004, Candès, Romberg, & Tao [77] and Donoho [136] released papers on perfect signal reconstruction from samples that stand in violation of Shannon's classical sampling theorem. These defiant signals are assumed sparse inherently or under some sparsifying affine transformation. Essentially, they proposed *sparse sampling theorems* asserting average sample rate independent of signal bandwidth and less than Shannon's rate.

MINIMUM SAMPLING RATE:

- OF  $\Omega$ -BANDLIMITED SIGNAL:  $2\Omega$  ([316, §3.2] Shannon)
- OF  $k$ -SPARSE LENGTH- $n$  SIGNAL:  $k \log_2(1+n/k)$  (Figure 113 Candès/Donoho)

Certainly, much was already known about nonuniform or random sampling [37] [288] and about subsampling or *multiprate systems* [99] [398]. Vetterli, Marziliano, & Blu [407] had congealed a theory of noiseless signal reconstruction, in May 2001, from samples that violate the Shannon rate. [427, *Sampling Sparsity*] They anticipated the sparsifying transform by recognizing: it is the *innovation* (onset) of functions constituting a (not necessarily bandlimited) signal that determines minimum sampling rate for perfect reconstruction. Average onset (sparsity), Vetterli *et alii* call, the *rate of innovation*. Vector inner-products that Candès/Donoho call *samples* or *measurements*, Vetterli calls *projections*. From those projections Vetterli demonstrates reconstruction (by digital signal processing and “root finding”) of a Dirac comb, the very same prototypical signal from which Candès probabilistically derives minimum sampling rate [*Compressive Sampling and Frontiers in Signal Processing*, University of Minnesota, June 6, 2007]. Combining their terminology, we paraphrase a sparse sampling theorem:

- Minimum sampling rate, asserted by Candès/Donoho,  $\propto$  Vetterli's rate of innovation (a.k.a: *information rate, degrees of freedom* [*ibidem*, June 5, 2007]).

What distinguishes these researchers are their methods of reconstruction.

Properties of the 1-norm were also well understood by June 2004 finding application in *deconvolution* of linear systems [91], constrained *linear regression (Lasso)* [388] [355], and *basis pursuit* [85] [240]. But never before had there been a formalized and rigorous sense that perfect reconstruction were possible by convex optimization of 1-norm when information lost in a subsampling process became nonrecoverable by classical methods.

Donoho named this discovery *compressed sensing* to describe a nonadaptive perfect reconstruction method by means of linear programming. By the time Candès' and Donoho's landmark papers were finally published by IEEE in 2006, compressed sensing was old news that had spawned intense research which still persists; notably, from prominent members of the *wavelet* community.

Reconstruction of the *Shepp-Logan phantom* (Figure 120), from a severely aliased image (Figure 122) obtained by Magnetic Resonance Imaging (MRI), was the *impetus* driving Candès' quest for a sparse sampling theorem. He realized that line segments appearing in the aliased image were regions of *high total variation*. There is great motivation, in the medical community, to apply compressed sensing to MRI because it translates to reduced scan-time which brings great technological and physiological benefits. MRI is now about 35 years old, beginning in 1973 with Nobel laureate Paul Lauterbur from Stony Brook USA. There has been much progress in MRI and compressed sensing since 2004, but there have also been indications of 1-norm abandonment (indigenous to reconstruction by compressed sensing) in favor of criteria closer to 0-norm because of a correspondingly smaller number of measurements required to accurately reconstruct a sparse signal:<sup>4.56</sup>

5481 complex samples (22 radial lines,  $\approx 256$  complex samples per) were required in June 2004 to reconstruct a noiseless  $256 \times 256$ -pixel Shepp-Logan phantom by 1-norm minimization of an image-gradient integral estimate called *total variation*; *id est*, 8.4% subsampling of 65536 data. [77, §1.1] [76, §3.2] It was soon discovered that reconstruction of the Shepp-Logan phantom were possible with only 2521 complex samples (10 radial lines, Figure 121); 3.8% subsampled data input to a (nonconvex)  $\frac{1}{2}$ -norm total-variation minimization. [83, §IIIA] The closer to 0-norm, the fewer the samples required for perfect reconstruction.

Passage of a few years witnessed an algorithmic speedup and dramatic reduction in minimum number of samples required for perfect reconstruction of the noiseless Shepp-Logan phantom. But minimization of total variation is ideally suited to recovery of any piecewise-constant image, like a phantom, because gradient of such images is highly sparse by design.

There is no inherent characteristic of real-life MRI images that would make reasonable an expectation of sparse gradient. Sparsification of a discrete image-gradient tends to preserve edges. Then minimization of total variation seeks an image having fewest edges. There is no deeper theoretical foundation than that. When applied to human brain scan or angiogram, with as much as 20% of  $256 \times 256$  Fourier samples, we have observed<sup>4.57</sup> a 30dB image/reconstruction-error ratio<sup>4.58</sup> barrier that seems impenetrable by the total-variation objective. Total-variation minimization has met with moderate success, in retrospect, only because some medical images are moderately piecewise-constant signals. One simply hopes a reconstruction, that is in some sense equal to a known subset of samples and whose gradient is most sparse, is that unique image we seek.<sup>4.59</sup>

---

<sup>4.56</sup>Efficient techniques continually emerge urging 1-norm criteria abandonment; [88] [397] [396, §IID] *e.g.*, five techniques for compressed sensing are compared in [38] demonstrating that 1-norm performance limits for cardinality minimization can be reliably exceeded.

<sup>4.57</sup>Experiments with real-life images were performed by Christine Law at Lucas Center for Imaging, Stanford University.

<sup>4.58</sup>Noise considered here is due only to the reconstruction process itself; *id est*, noise in excess of that produced by the best reconstruction of an image from a complete set of samples in the sense of Shannon. At less than 30dB image/error, artifacts generally remain visible to the naked eye. We estimate that about 50dB is required to eliminate noticeable distortion in a visual A/B comparison.

<sup>4.59</sup>In vascular radiology, diagnoses are almost exclusively based on morphology of vessels and, in particular, presence of stenoses. There is a compelling argument for total-variation reconstruction of magnetic resonance angiogram because it helps isolate structures of particular interest.

The total-variation objective, operating on an image, is expressible as norm of a linear transformation (934). It is natural to ask whether there exist other sparsifying transforms that might break the real-life 30dB barrier (any sampling pattern @20%  $256 \times 256$  data) in MRI. There has been much research into application of wavelets, discrete cosine transform (DCT), randomized orthogonal bases, splines, *etcetera*, but with suspiciously little focus on objective measures like image/error or illustration of difference images; the predominant basis of comparison instead being subjectively visual (Duensing & Huang, ISMRM Toronto 2008).<sup>4.60</sup> Despite choice of transform, there seems yet to have been a breakthrough of the 30dB barrier. Application of compressed sensing to MRI, therefore, remains fertile in 2008 for continued research.

### regularized form of compressed sensing in imaging

We now repeat Candès' image reconstruction experiment from 2004 which led to discovery of sparse sampling theorems. [77, §1.2] But we achieve perfect reconstruction with an algorithm based on vanishing gradient of a compressed sensing problem's regularization, which is computationally efficient. Our *contraction method* (p.305) is fast also because matrix multiplications are replaced by fast Fourier transform, and number of constraints is cut in half by sampling symmetrically. Convex iteration for cardinality minimization (§4.6) is incorporated which allows perfect reconstruction of a phantom at 4.1% subsampling rate; 50% Candès' rate. By making neighboring-pixel selection adaptive, convex iteration reduces discrete image-gradient sparsity of the Shepp-Logan phantom to 1.9%; 33% lower than previously reported.

We demonstrate application of discrete image-gradient sparsification to the  $n \times n = 256 \times 256$  Shepp-Logan phantom, simulating idealized acquisition of MRI data by radial sampling in the Fourier domain (Figure 121).<sup>4.61</sup> Define a Nyquist-centric *discrete Fourier transform* (DFT) matrix

$$F \triangleq \begin{bmatrix} 1 & 1 & 1 & 1 & \cdots & 1 \\ 1 & e^{-j2\pi/n} & e^{-j4\pi/n} & e^{-j6\pi/n} & \cdots & e^{-j(n-1)2\pi/n} \\ 1 & e^{-j4\pi/n} & e^{-j8\pi/n} & e^{-j12\pi/n} & \cdots & e^{-j(n-1)4\pi/n} \\ 1 & e^{-j6\pi/n} & e^{-j12\pi/n} & e^{-j18\pi/n} & \cdots & e^{-j(n-1)6\pi/n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & e^{-j(n-1)2\pi/n} & e^{-j(n-1)4\pi/n} & e^{-j(n-1)6\pi/n} & \cdots & e^{-j(n-1)^2 2\pi/n} \end{bmatrix} \frac{1}{\sqrt{n}} \in \mathbb{C}^{n \times n} \quad (916)$$

a symmetric (nonHermitian) unitary matrix characterized

$$\begin{aligned} F &= F^T \\ F^{-1} &= F^H \end{aligned} \quad (917)$$

Denoting an unknown image  $\mathcal{U} \in \mathbb{R}^{n \times n}$ , its two-dimensional discrete Fourier transform  $\mathfrak{F}$  is

$$\mathfrak{F}(\mathcal{U}) \triangleq F \mathcal{U} F^H \quad (918)$$

hence the inverse discrete transform

$$\mathcal{U} = F^H \mathfrak{F}(\mathcal{U}) F^H \quad (919)$$

---

<sup>4.60</sup> I have never calculated the PSNR of these reconstructed images [of Barbara]. —Jean-Luc Starck  
The sparsity of the image is the percentage of transform coefficients sufficient for diagnostic-quality reconstruction. Of course the term “diagnostic quality” is subjective. . . . I have yet to see an “objective” measure of image quality. Difference images, in my experience, definitely do not tell the whole story. Often I would show people some of my results and get mixed responses, but when I add artificial Gaussian noise to an image, often people say that it looks better. —Michael Lustig

<sup>4.61</sup>  $k$ -space is conventional acquisition terminology indicating domain of the continuous raw data provided by an MRI machine. An image is reconstructed by inverse discrete Fourier transform of that data interpolated on a Cartesian grid in two dimensions.

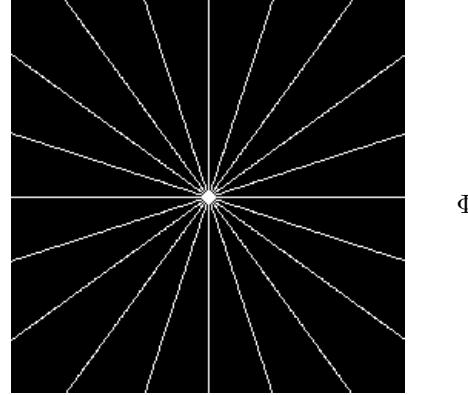


Figure 121: MRI radial sampling pattern, in DC-centric Fourier domain, representing 4.1% (10 lines) subsampled data. Only half of these complex samples, in any halfspace about the origin in theory, need be acquired for a real image because of conjugate symmetry. Due to MRI machine imperfections, samples are generally taken over full extent of each radial line segment. MRI acquisition time is proportional to number of lines.

From §A.1.1 no.33 we have a vectorized two-dimensional DFT via Kronecker product  $\otimes$

$$\text{vec } \mathfrak{F}(\mathcal{U}) \triangleq (F \otimes F) \text{vec } \mathcal{U} \quad (920)$$

and from (919) its inverse [190, p.24]

$$\text{vec } \mathcal{U} = (F^H \otimes F^H)(F \otimes F) \text{vec } \mathcal{U} = (F^H F \otimes F^H F) \text{vec } \mathcal{U} \quad (921)$$

Idealized radial sampling in the Fourier domain can be simulated by Hadamard product  $\circ$  with a binary mask  $\Phi \in \mathbb{R}^{n \times n}$  whose nonzero entries could, for example, correspond with the radial line segments in Figure 121. To make the mask Nyquist-centric, like DFT matrix  $F$ , define a circulant [192] symmetric permutation matrix<sup>4.62</sup>

$$\Theta \triangleq \begin{bmatrix} \mathbf{0} & I \\ I & \mathbf{0} \end{bmatrix} \in \mathbb{S}^n \quad (922)$$

Then given subsampled Fourier domain (MRI  $\mathbb{k}$ -space) measurements in incomplete  $K \in \mathbb{C}^{n \times n}$ , we might constrain  $\mathfrak{F}(\mathcal{U})$  thus:

$$\Theta \Phi \Theta \circ F \mathcal{U} F = K \quad (923)$$

and in vector form, (42) (1963)

$$\delta(\text{vec } \Theta \Phi \Theta)(F \otimes F) \text{vec } \mathcal{U} = \text{vec } K \quad (924)$$

Because measurements  $K$  are complex, there are actually twice the number of equality constraints as there are measurements.

We can cut that number of constraints in half via vertical and horizontal mask  $\Phi$  symmetry which forces the imaginary inverse transform to  $\mathbf{0}$ : The inverse subsampled transform in matrix form is

$$F^H (\Theta \Phi \Theta \circ F \mathcal{U} F) F^H = F^H K F^H \quad (925)$$

---

<sup>4.62</sup>MATLAB `fftshift()`

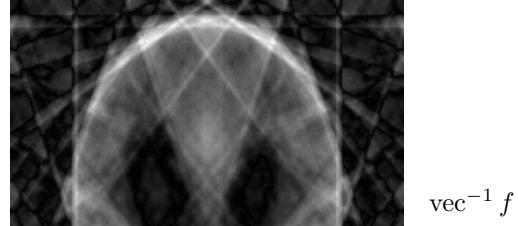


Figure 122: Aliasing of Shepp-Logan phantom in Figure 120 resulting from  $k$ -space subsampling pattern in Figure 121. This image is real because binary mask  $\Phi$  is vertically and horizontally symmetric. It is remarkable that the phantom can be reconstructed, by convex iteration, given only  $\mathcal{U}_0 = \text{vec}^{-1} f$ .

and in vector form

$$(F^H \otimes F^H) \delta(\text{vec } \Theta \Phi \Theta) (F \otimes F) \text{vec } \mathcal{U} = (F^H \otimes F^H) \text{vec } K \quad (926)$$

later abbreviated

$$P \text{vec } \mathcal{U} = f \quad (927)$$

where

$$P \triangleq (F^H \otimes F^H) \delta(\text{vec } \Theta \Phi \Theta) (F \otimes F) \in \mathbb{C}^{n^2 \times n^2} \quad (928)$$

Because of idempotence  $P = P^2$ ,  $P$  is a projection matrix. Because of its Hermitian symmetry [190, p.24]

$$P = (F^H \otimes F^H) \delta(\text{vec } \Theta \Phi \Theta) (F \otimes F) = (F \otimes F)^H \delta(\text{vec } \Theta \Phi \Theta) (F^H \otimes F^H)^H = P^H \quad (929)$$

$P$  is an orthogonal projector.<sup>4.63</sup>  $P \text{vec } \mathcal{U}$  is real when  $P$  is real; *id est*, when for positive even integer  $n$

$$\Phi = \begin{bmatrix} \Phi_{11} & \Phi(1, 2:n) \Xi \\ \Xi \Phi(2:n, 1) & \Xi \Phi(2:n, 2:n) \Xi \end{bmatrix} \in \mathbb{R}^{n \times n} \quad (930)$$

where  $\Xi \in \mathbb{S}^{n-1}$  is the order-reversing permutation matrix (1900). In words, this necessary and sufficient condition on  $\Phi$  (for a real inverse subsampled transform [316, p.53]) demands vertical symmetry about row  $\frac{n}{2}+1$  and horizontal symmetry<sup>4.64</sup> about column  $\frac{n}{2}+1$ .

Define

$$\Delta \triangleq \begin{bmatrix} 1 & 0 & & & \mathbf{0} \\ -1 & 1 & 0 & & \\ & -1 & 1 & \ddots & \\ & & \ddots & \ddots & \ddots \\ & & & \ddots & 1 & 0 \\ \mathbf{0}^T & & & & -1 & 1 \end{bmatrix} \in \mathbb{R}^{n \times n} \quad (931)$$

<sup>4.63</sup> (928) is a diagonalization of matrix  $P$  whose binary eigenvalues are  $\delta(\text{vec } \Theta \Phi \Theta)$  while the corresponding eigenvectors constitute the columns of unitary matrix  $F^H \otimes F^H$ .

<sup>4.64</sup>This condition on  $\Phi$  applies to both DC- and Nyquist-centric DFT matrices.

Express an image-gradient estimate

$$\nabla \mathcal{U} \triangleq \begin{bmatrix} \mathcal{U} \Delta \\ \mathcal{U} \Delta^T \\ \Delta \mathcal{U} \\ \Delta^T \mathcal{U} \end{bmatrix} \in \mathbb{R}^{4n \times n} \quad (932)$$

that is a simple first-order difference of neighboring pixels (Figure 123) to the right, left, above, and below.<sup>4.65</sup> By §A.1.1 no.33, its vectorization: for  $\Psi_i \in \mathbb{R}^{n^2 \times n^2}$

$$\text{vec } \nabla \mathcal{U} = \begin{bmatrix} \Delta^T \otimes I \\ \Delta \otimes I \\ I \otimes \Delta \\ I \otimes \Delta^T \end{bmatrix} \text{vec } \mathcal{U} \triangleq \begin{bmatrix} \Psi_1 \\ \Psi_1^T \\ \Psi_2 \\ \Psi_2^T \end{bmatrix} \text{vec } \mathcal{U} \triangleq \Psi \text{vec } \mathcal{U} \in \mathbb{R}^{4n^2} \quad (933)$$

where  $\Psi \in \mathbb{R}^{4n^2 \times n^2}$ . A total-variation minimization for reconstructing MRI image  $\mathcal{U}$ , that is known suboptimal [234] [78], may be concisely posed

$$\begin{aligned} & \underset{\mathcal{U}}{\text{minimize}} \quad \|\Psi \text{vec } \mathcal{U}\|_1 \\ & \text{subject to} \quad P \text{vec } \mathcal{U} = f \end{aligned} \quad (934)$$

where

$$f = (F^H \otimes F^H) \text{vec } K \in \mathbb{C}^{n^2} \quad (935)$$

is the known inverse subsampled Fourier data (a vectorized aliased image, Figure 122), and where a norm of discrete image-gradient  $\nabla \mathcal{U}$  is equivalently expressed as norm of a linear transformation  $\Psi \text{vec } \mathcal{U}$ .

Although this simple problem statement (934) is equivalent to a linear program (§3.2), its numerical solution is beyond the capability of even the most highly regarded of contemporary commercial solvers.<sup>4.66</sup> Our recourse is to recast the problem in regularized form and write customized code to solve it:

$$\begin{aligned} & \underset{\mathcal{U}}{\text{minimize}} \quad \langle |\Psi \text{vec } \mathcal{U}|, y \rangle \\ & \text{subject to} \quad P \text{vec } \mathcal{U} = f \\ & \equiv \\ & \underset{\mathcal{U}}{\text{minimize}} \quad \langle |\Psi \text{vec } \mathcal{U}|, y \rangle + \frac{1}{2} \lambda \|P \text{vec } \mathcal{U} - f\|_2^2 \end{aligned} \quad \begin{aligned} & (a) \\ & (b) \end{aligned} \quad (936)$$

where multiobjective parameter  $\lambda \in \mathbb{R}_+$  is quite large ( $\lambda \approx 1E8$ ) so as to enforce the equality constraint:  $P \text{vec } \mathcal{U} - f = \mathbf{0} \Leftrightarrow \|P \text{vec } \mathcal{U} - f\|_2^2 = 0$  (§A.7.1). We introduce a direction vector  $y \in \mathbb{R}_+^{4n^2}$  as part of a convex iteration (§4.6.3) to overcome that known suboptimal minimization of discrete image-gradient cardinality: *id est*, there exists a vector  $y^*$  with entries  $y_i^* \in \{0, 1\}$  such that

$$\begin{aligned} & \underset{\mathcal{U}}{\text{minimize}} \quad \|\Psi \text{vec } \mathcal{U}\|_0 \\ & \text{subject to} \quad P \text{vec } \mathcal{U} = f \end{aligned} \quad \equiv \quad \underset{\mathcal{U}}{\text{minimize}} \quad \langle |\Psi \text{vec } \mathcal{U}|, y^* \rangle + \frac{1}{2} \lambda \|P \text{vec } \mathcal{U} - f\|_2^2 \quad (937)$$

Existence of such a  $y^*$ , complementary to an optimal vector  $\Psi \text{vec } \mathcal{U}^*$ , is obvious by definition of global optimality  $\langle |\Psi \text{vec } \mathcal{U}^*|, y^* \rangle = 0$  (837) under which a cardinality- $c$  optimal objective  $\|\Psi \text{vec } \mathcal{U}^*\|_0$  is assumed to exist.

<sup>4.65</sup>There is significant improvement in reconstruction quality by augmentation of a nominally two-point discrete image-gradient estimate to four points per pixel by inclusion of two polar directions. Improvement is due to centering; symmetry of discrete differences about a central pixel. We find small improvement on real-life images,  $\approx 1\text{dB}$  empirically, by further augmentation with diagonally adjacent pixel differences.

<sup>4.66</sup>for images as small as  $128 \times 128$  pixels. Obstacle to numerical solution is not a computer resource: *e.g.*, execution time, memory. The obstacle is, in fact, inadequate numerical precision. Even when all dependent equality constraints are manually removed, the best commercial solvers fail simply because computer numerics become nonsense; *id est*, numerical errors enter significant digits and the algorithm exits prematurely, loops indefinitely, or produces an infeasible solution.

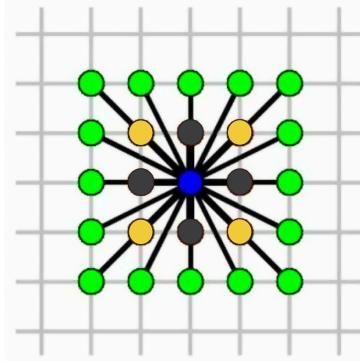


Figure 123: Neighboring-pixel stencil [397] for image-gradient estimation on Cartesian grid. Implementation selects adaptively from darkest four  $\bullet$  about central. Continuous image-gradient from two pixels holds only in a limit. For discrete differences, better practical estimates are obtained when centered.

Because (936b) is an unconstrained convex problem, a zero objective function gradient is necessary and sufficient for optimality (§2.13.3); *id est*, (§D.2.1)

$$\Psi^T \delta(y) \operatorname{sgn}(\Psi \operatorname{vec} \mathcal{U}) + \lambda P^H (P \operatorname{vec} \mathcal{U} - f) = \mathbf{0} \quad (938)$$

Because of  $P$  idempotence and Hermitian symmetry and  $\operatorname{sgn}()$  definition (p.633), this is equivalent to

$$\lim_{\epsilon \rightarrow 0} (\Psi^T \delta(y) \delta(|\Psi \operatorname{vec} \mathcal{U}| + \epsilon \mathbf{1})^{-1} \Psi + \lambda P) \operatorname{vec} \mathcal{U} = \lambda P f \quad (939)$$

where small positive constant  $\epsilon \in \mathbb{R}_+$  has been introduced for invertibility. Speaking more analytically, introduction of  $\epsilon$  serves to uniquely define the objective's gradient everywhere in the function domain; *id est*, it transforms absolute value in (936b) from a function differentiable almost everywhere into a differentiable function. An example of such a transformation in one dimension is illustrated in Figure 124. When small enough for practical purposes<sup>4.67</sup> ( $\epsilon \approx 1E-3$ ), we may ignore the limiting operation. Then the mapping, for  $0 \preceq y \preceq \mathbf{1}$

$$\operatorname{vec} \mathcal{U}_{t+1} = (\Psi^T \delta(y) \delta(|\Psi \operatorname{vec} \mathcal{U}_t| + \epsilon \mathbf{1})^{-1} \Psi + \lambda P)^{-1} \lambda P f \quad (940)$$

is a *contraction* in  $\mathcal{U}_t$  that can be solved recursively in  $t$  for its unique *fixed point*; *id est*, until  $\mathcal{U}_{t+1} \rightarrow \mathcal{U}_t$ . [254, p.300] [229, p.155] Calculating this inversion directly is not possible for large matrices on contemporary computers because of numerical precision, so instead we apply the *conjugate gradient* method of solution to

$$(\Psi^T \delta(y) \delta(|\Psi \operatorname{vec} \mathcal{U}_t| + \epsilon \mathbf{1})^{-1} \Psi + \lambda P) \operatorname{vec} \mathcal{U}_{t+1} = \lambda P f \quad (941)$$

which is linear in  $\mathcal{U}_{t+1}$  at each recursion in the MATLAB program [412].<sup>4.68</sup>

<sup>4.67</sup>We are looking for at least 50dB image/error ratio from only 4.1% subsampled data (10 radial lines in  $\mathbb{k}$ -space). With this setting of  $\epsilon$ , we actually attain in excess of 100dB from a simple MATLAB program in about a minute on a 2006 vintage laptop Core 2 CPU (Intel T7400@2.16GHz, 666MHz FSB). By trading execution time and treating discrete image-gradient cardinality as a known quantity for this phantom, over 160dB is achievable.

<sup>4.68</sup>Conjugate gradient method requires positive definiteness. [173, §4.8.3.2]

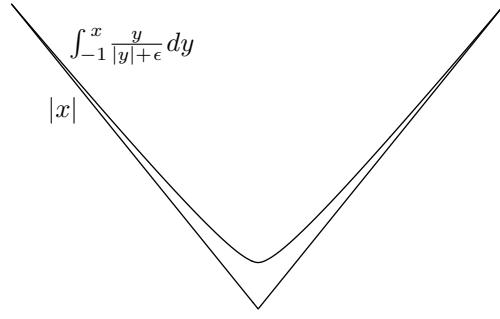


Figure 124: Real absolute value function  $f_2(x) = |x|$  on  $x \in [-1, 1]$  from Figure 72b superimposed upon integral of its derivative at  $\epsilon = 0.05$  which smooths objective function.

Observe that  $P$  (928), in the equality constraint from problem (936a), is not a wide matrix.<sup>4.69</sup> Although number of Fourier samples taken is equal to the number of nonzero entries in binary mask  $\Phi$ , matrix  $P$  is square but never actually formed during computation. Rather, a two-dimensional fast Fourier transform of  $\mathcal{U}$  is computed followed by masking with  $\Theta\Phi\Theta$  and then an inverse fast Fourier transform. This technique significantly reduces memory requirements and, together with contraction method of solution, is the principal reason for relatively fast computation.

### convex iteration

By *convex iteration* we mean alternation of solution to (936a) and (942) until convergence. Direction vector  $y$  is initialized to  $\mathbf{1}$  until the first fixed point is found; which means, the contraction recursion begins calculating a (1-norm) solution  $\mathcal{U}^*$  to (934) via problem (936b). Once  $\mathcal{U}^*$  is found, vector  $y$  is updated according to an estimate of discrete image-gradient cardinality  $c$ : Sum of the  $4n^2 - c$  smallest entries of  $|\Psi \text{vec } \mathcal{U}^*| \in \mathbb{R}^{4n^2}$  is the optimal objective value from a linear program, for  $0 \leq c \leq 4n^2 - 1$  (530)

$$\begin{aligned} \sum_{i=c+1}^{4n^2} \pi(|\Psi \text{vec } \mathcal{U}^*|)_i &= \underset{\substack{y \in \mathbb{R}^{4n^2} \\ \text{subject to}}} {\underset{\substack{0 \preceq y \preceq \mathbf{1} \\ y^T \mathbf{1} = 4n^2 - c}}{\text{minimize}}} \quad \langle |\Psi \text{vec } \mathcal{U}^*|, y \rangle \end{aligned} \quad (942)$$

where  $\pi$  is the nonlinear permutation-operator sorting its vector argument into nonincreasing order. An *optimal solution*  $y$  to (942), that is an extreme point of its feasible set, is known in closed form: it has 1 in each entry corresponding to the  $4n^2 - c$  smallest entries of  $|\Psi \text{vec } \mathcal{U}^*|$  and has 0 elsewhere. – p.273 Updated image  $\mathcal{U}^*$  is assigned to  $\mathcal{U}_t$ , the contraction is recomputed solving (936b), direction vector  $y$  is updated again, and so on until convergence which is guaranteed by virtue of a monotonically nonincreasing real sequence of objective values in (936a) and (942).

There are two features that distinguish problem formulation (936b) and our particular implementation of it [412, MATLAB code]:

---

<sup>4.69</sup>Wide is typical of compressed sensing problems; e.g, [76] [83].

- 1) An image-gradient estimate may engage any combination of four adjacent pixels. In other words, the algorithm is not locked into a four-point gradient estimate (Figure 123); number of points constituting an estimate is directly determined by direction vector  $y$ .<sup>4.70</sup> Indeed, we find only  $c = 5092$  zero entries in  $y^*$  for the Shepp-Logan phantom; meaning, discrete image-gradient sparsity is actually closer to 1.9% than the 3% reported elsewhere; e.g., [396, §IIB].
- 2) Numerical precision of the fixed point of contraction (940) ( $\approx 1E-2$  for perfect reconstruction @ -103dB error) is a parameter to the implementation; meaning, direction vector  $y$  is updated after contraction begins but prior to its culmination. Impact of this idiosyncrasy tends toward simultaneous optimization in variables  $\mathcal{U}$  and  $y$  while insuring  $y$  settles on a boundary point of its feasible set (nonnegative hypercube slice) in (942) at every iteration; for only a boundary point<sup>4.71</sup> can yield the sum of smallest entries in  $|\Psi \text{vec } \mathcal{U}^*|$ .

Perfect reconstruction of the Shepp-Logan phantom (at 103dB image/error) is achieved in a MATLAB minute with 4.1% subsampled data (2671 complex samples); well below an 11% least lower bound predicted by the sparse sampling theorem. Because reconstruction approaches optimal solution to a 0-norm problem, minimum number of Fourier-domain samples is bounded below by cardinality of discrete image-gradient at 1.9%.  $\square$

#### 4.7.0.0.14 Exercise. *Contraction operator.*

Determine conditions on  $\lambda$  and  $\epsilon$  under which  $\Psi^T \delta(y) \delta(|\Psi \text{vec } \mathcal{U}_t| + \epsilon \mathbf{1})^{-1} \Psi + \lambda P$  from (941) is positive definite and (940) is a contraction.  $\blacktriangledown$

#### 4.7.0.0.15 Example. *Eternity II.*

A tessellation puzzle game, playable by children, commenced world-wide in July 2007; introduced in London by Christopher Walter Monckton, 3<sup>rd</sup> Viscount Monckton of Brenchley. Called Eternity II, its name derives from an estimate of time that would pass while trying all allowable tilings of puzzle pieces before obtaining a complete solution. By the end of 2008, a complete solution had not yet been found although a \$10,000 USD prize was awarded for a high score 467 (out of  $480 = 2\sqrt{M}(\sqrt{M} - 1)$ ) obtained by heuristic methods.<sup>4.72</sup> No prize was awarded for 2009 and 2010. Game-rules state that a \$2M prize would be awarded to the first person who completely solves the puzzle before December 31, 2010, but the prize went unclaimed and solution remains yet to be found.

The full game comprises  $M = 256$  square pieces and  $16 \times 16$  gridded board (Figure 126) whose complete tessellation is considered NP-hard.<sup>4.73</sup> [382] [119] A player may tile, retile, and rotate pieces, indexed 1 through 256, in any order face-up on the square board. Pieces are immutable in the sense that each is characterized by four colors (and their uniquely associated British symbols), one at each edge, which are not necessarily the same per piece or from piece to piece; *id est*, different pieces may or may not have some edge-colors in common. There are  $L = 22$  distinct edge-colors plus a solid grey. The object of the game is to completely tile the board with pieces whose touching edges have identical color. Boundary of the board must be colored grey.

<sup>4.70</sup>This adaptive gradient was not contrived. It is an artifact of the convex iteration method for minimal cardinality solution; in this case, cardinality minimization of a discrete image-gradient.

<sup>4.71</sup>Simultaneous optimization of these two variables  $\mathcal{U}$  and  $y$  should never be a pinnacle of aspiration; for then, optimal  $y$  might not attain a boundary point.

<sup>4.72</sup>That score means all but a few of the 256 pieces had been placed successfully (including the mandatory piece). Although distance between 467 to 480 is relatively small, there is apparently vast distance to a solution because no complete solution followed in 2009.

<sup>4.73</sup>Even so, combinatorial-intensity brute-force backtracking methods can solve similar puzzles in minutes given  $M = 196$  pieces on a  $14 \times 14$  test board; as demonstrated by [Yannick Kirschhofer](#). There is a steep rise in level of difficulty going to a  $15 \times 15$  board.

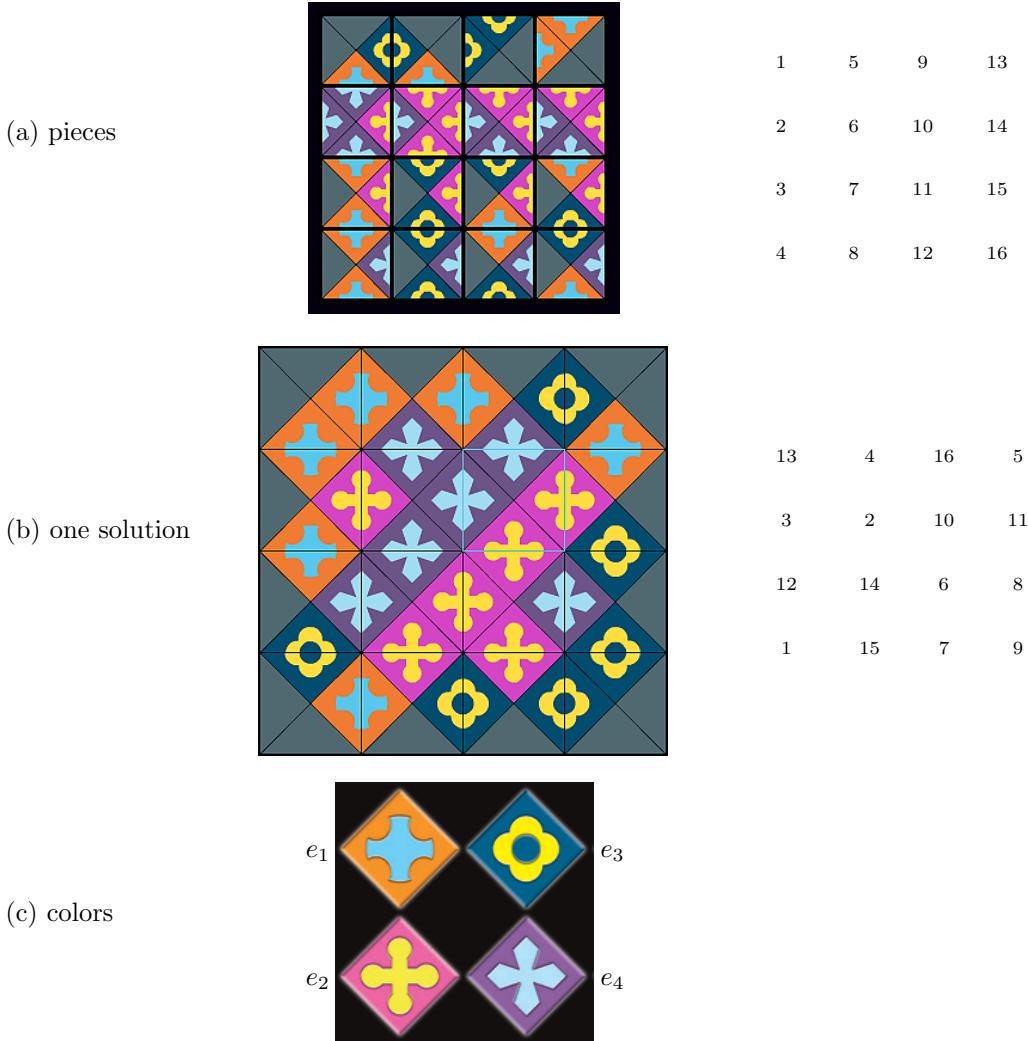


Figure 125: *Eternity* II is a board game in the puzzle genre. **(a)** Shown are all of the 16 puzzle pieces (indexed as in the tableau alongside) from a scaled-down computerized demonstration game version on the [TOMY website](#). Puzzle pieces are square and partitioned into four colors (with associated symbols). Pieces may be moved, removed, and rotated at random on a  $4 \times 4$  board. **(b)** Illustrated is one complete solution to this puzzle whose solution is not unique. The piece, whose border is lightly outlined, was placed last in this realization. There is no mandatory piece placement, as for the full game, except the grey board-boundary. Solution time for a human is typically on the order of a minute. **(c)** This puzzle has four colors, indexed 1 through 4; grey corresponds to 0.

### full-game rules

- 1) Any puzzle piece may be rotated face-up in quadrature and placed or replaced on the square board.
- 2) Only one piece may occupy any particular cell on the board.
- 3) All adjacent pieces must match in color (and symbol) at their touching edges.
- 4) Solid grey edges must appear all along the board's boundary.
- 5) One mandatory piece (numbered 139) must have a predetermined rotation in a predetermined cell (number 121) on the board (Figure 126).
- 6) The board must be tiled completely (*covered*).

A scaled-down demonstration version of the game is illustrated in Figure 125. Differences between the full game (Figure 126) and scaled-down game are: number of edge-colors  $L$  (22 *versus* 4, ignoring solid grey), number of pieces  $M$  (256 *versus* 16), and a single mandatory piece placement interior to the board for the full game. The scaled-down game has four distinct edge-colors, plus a solid grey, whose coding is illustrated in Figure 125c.

- For the full game board, there are  $L = 22$  distinct edge-colors and  $M = 256$  puzzle pieces with board-dimension  $\sqrt{M} \times \sqrt{M} = 16 \times 16$ .
- For the scaled-down demonstration game board, there are  $L = 4$  distinct edge-colors and  $M = 16$  puzzle pieces with board-dimension  $\sqrt{M} \times \sqrt{M} = 4 \times 4$ .

### Euclidean distance intractability

If each square puzzle piece were characterized by four points in quadrature, one point representing board coordinates and color per edge, then Euclidean distance geometry would be suitable for solving this puzzle. Since all interpoint distances per piece are known, this game may be regarded as a Euclidean distance matrix completion problem<sup>4.74</sup> in  $\text{EDM}^{4M}$ . Because distance information provides for reconstruction of point position to within an isometry (§5.5), piece translation and rotation are isometric transformations that abide by rules of the game.<sup>4.75</sup> Convex constraints can be devised to prevent puzzle-piece reflection and to quantize rotation such that piece-edges stay aligned with the board boundary. (§5.5.2.0.1)

But manipulating such a large EDM is too numerically difficult for contemporary general-purpose semidefinite program (SDP) solvers which incorporate interior-point methods; indeed, they are hard-pressed to find a solution for variable matrices of dimension as small as 100. Our challenge, therefore, is to express this game's rules as constraints in a convex and numerically tractable way so as to find one solution from a googol of possible combinations.<sup>4.76</sup>

---

<sup>4.74</sup>(§6.7) Were edge-points ordered sequentially with piece number, then this EDM would have a block-diagonal structure of known entries.

<sup>4.75</sup>Translation occurs when a piece moves on the board in Figure 126, rotation occurs when colors are aligned with an adjacent piece.

<sup>4.76</sup>Oliver Riordan asserts that at least one solution exists; I suspect there is only one solution although Monckton insists they number in the thousands. Ignoring board-boundary constraints and the full game's single mandatory piece placement, a loose upper bound on number of combinations is  $M! 4^M = 256! 4^{256}$ . That number gets further loosened:  $150638!/(256!(150638 - 256)!)$  after presolving Eternity II (970).

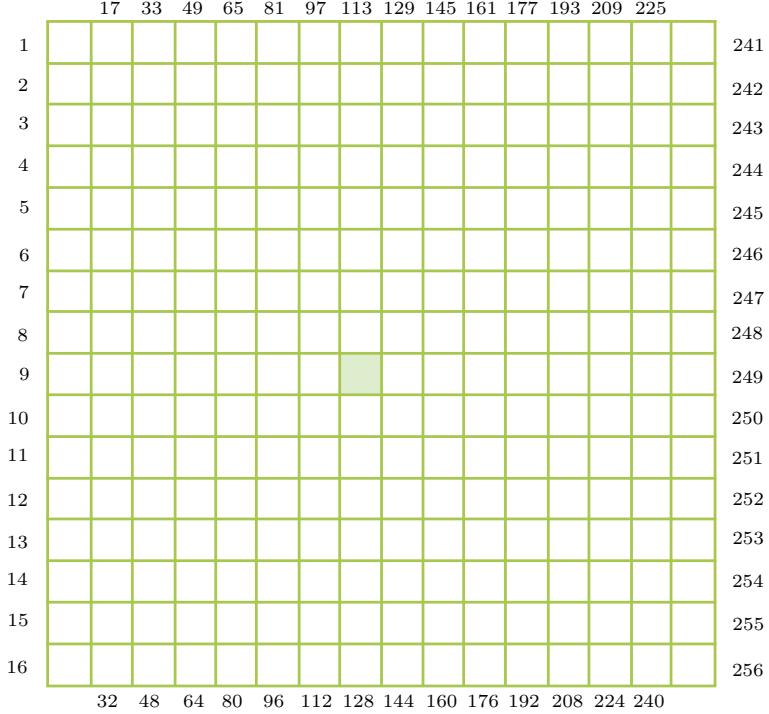


Figure 126: *Eternity II* full-game board ( $16 \times 16$ ,  $M = 256$ ,  $L = 22$ ) illustrating boundary cell numbers. Grid facilitates piece placement within unit-square cell; one piece per cell. Cell 121 (shaded) holds mandatory puzzle-piece  $P_{139}$  designated by Monckton.

#### piece $P$ permutation $\Xi$ rotation $\Pi$ strategy

To each puzzle piece, from a given set of  $M$  pieces  $\{P_i, i=1 \dots M\}$ , assign an index  $i$  representing a unique piece-number. Each square piece is characterized by four given colors, in quadrature, corresponding to its four edges. Each color  $p_{ij} \in \mathbb{R}^L$  is represented by  $e_\ell \in \mathbb{R}^L$  an  $L$ -dimensional standard basis vector or  $\mathbf{0}$  if grey. These four edge-colors are represented in a  $4 \times L$ -dimensional matrix; one matrix per piece

$$P_i \triangleq \begin{bmatrix} p_{i1}^T \\ p_{i2}^T \\ p_{i3}^T \\ p_{i4}^T \end{bmatrix} \in \mathbb{R}^{4 \times L}, \quad i=1 \dots M \quad (943)$$

In other words, each distinct nongrey color is assigned a unique corresponding index  $\ell \in \{1 \dots L\}$  identifying a standard basis vector  $e_\ell \in \mathbb{R}^L$  (Figure 125c) that becomes a vector  $p_{ij} \in \{e_1 \dots e_L, \mathbf{0}\} \subset \mathbb{R}^L$  constituting matrix  $P_i$  representing a particular piece. Rows  $\{p_{ij}^T, j=1 \dots 4\}$  of  $P_i$  are ordered counterclockwise as in Figure 127. Color data is given in Figure 128 for the demonstration game board. Then matrix  $P_i$  describes the  $i^{\text{th}}$  piece, excepting its rotation and position on the board.

Our intent is to show how to vectorize the board, with respect to whole pieces, and then express Eternity II as a very hard combinatorial objective with linear constraints: All pieces are initially placed in order of their given index  $i$  assigned by Monckton. The vectorized game-board has initial state represented within a matrix

$$P_6 = \begin{bmatrix} p_{61}^T \\ p_{62}^T \\ p_{63}^T \\ p_{64}^T \end{bmatrix} \in \mathbb{R}^{4 \times L}$$

Figure 127: Demo-game piece  $P_6$  illustrating edge-color  $\bullet$   $p_{6j} \in \mathbb{R}^L$  counterclockwise ordering in  $j$  beginning from right. For all game boards, edge-color index  $j=1 \dots 4$ .

$$P \triangleq \begin{bmatrix} P_1 \\ \vdots \\ P_M \end{bmatrix} \in \mathbb{R}^{4M \times L} \quad (944)$$

enumerated in Figure 128 for the demonstration game. Moving pieces all at once about the square board corresponds to permuting pieces  $P_i$  on the vectorized board represented by matrix  $P$ , while rotating the  $i^{\text{th}}$  piece is equivalent to circularly shifting row indices of  $P_i$  (rowwise permutation). This permutation problem, as stated, is doubly combinatorial ( $M! 4^M$  combinations) because we must find a permutation of pieces ( $M!$ )

$$\Xi \in \mathbb{R}^{M \times M} \quad (945)$$

and quadrature rotation  $\Pi_i \in \mathbb{R}^{4 \times 4}$  of each individual piece ( $4^M$ ) that solve the puzzle;

$$(\Xi \otimes I_4) \Pi P = (\Xi \otimes I_4) \begin{bmatrix} \Pi_1 P_1 \\ \vdots \\ \Pi_M P_M \end{bmatrix} \in \mathbb{R}^{4M \times L} \quad (946)$$

where

$$\Pi_i \in \{\pi_1, \pi_2, \pi_3, \pi_4\} \triangleq \left\{ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \right\} \quad (947)$$

$$\Pi \triangleq \begin{bmatrix} \Pi_1 & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & & \Pi_M \end{bmatrix} \in \mathbb{R}^{4M \times 4M} \quad (948)$$

and where  $I_4 \triangleq I \in \mathbb{S}^4$  and  $\pi_1 = I_4$ . Initial game-board state  $P$  (944) corresponds to  $\Xi = I$  and  $\Pi_i = \pi_1 \forall i$ . Circulant [192] permutation matrices  $\{\pi_1, \pi_2, \pi_3, \pi_4\} \subset \mathbb{R}^{4 \times 4}$  correspond to clockwise piece-rotations  $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$ .

### piece edge adjacency $\Delta$

Rules of the game dictate that adjacent pieces on the square board have colors that match at their touching edges as in Figure 125b.<sup>4.77</sup> A complete match is therefore equivalent to demanding that a constraint, comprising numeric color differences between  $2\sqrt{M}(\sqrt{M}-1)$  touching edges, vanish. Because vectorized board layout is fixed and its cells are loaded or reloaded with pieces during play, locations of adjacent edges in  $\mathbb{R}^{4M \times L}$  (946) are known *a priori*. We need simply form differences between colors from adjacent edges of pieces loaded into those known locations. Each difference may be represented

<sup>4.77</sup>Piece adjacencies on the square board map linearly to the vectorized board, of course.

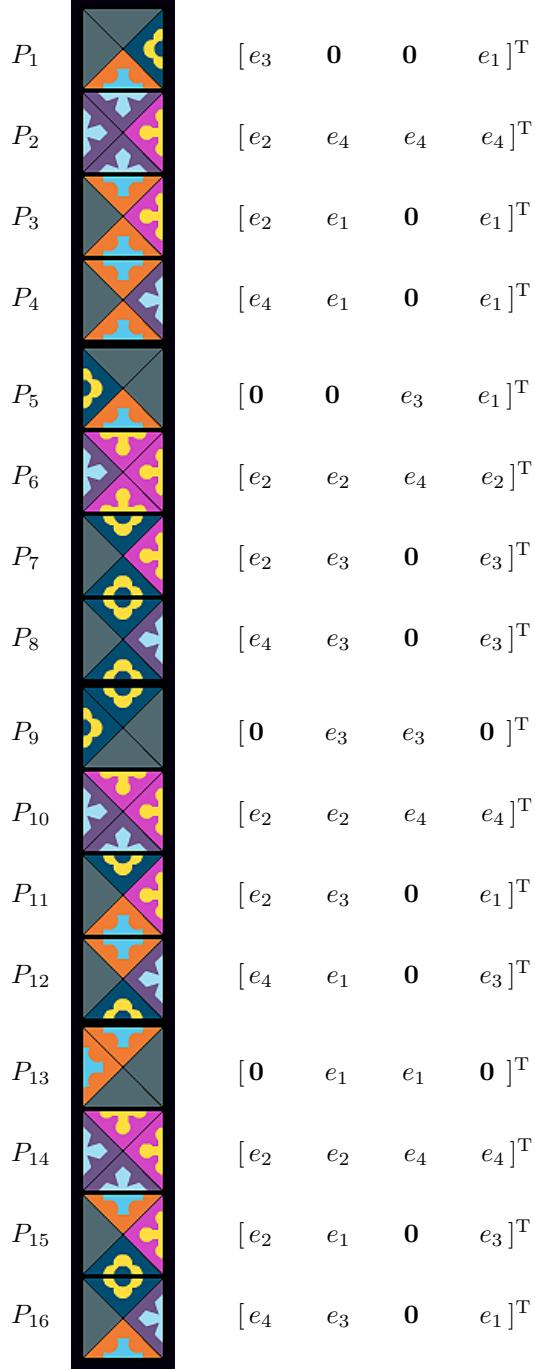


Figure 128: Vectorized demo-game board illustrating  $M=16$  matrices in  $\mathbb{R}^{4 \times L}$  describing initial state  $P \in \mathbb{R}^{4M \times L}$  of puzzle pieces; four colors per puzzle-piece (Figure 127),  $L=4$  colors total in game (Figure 125c). Standard basis vectors  $e_\ell$  in  $\mathbb{R}^L$  represent color so that color difference measurement remains unweighted.

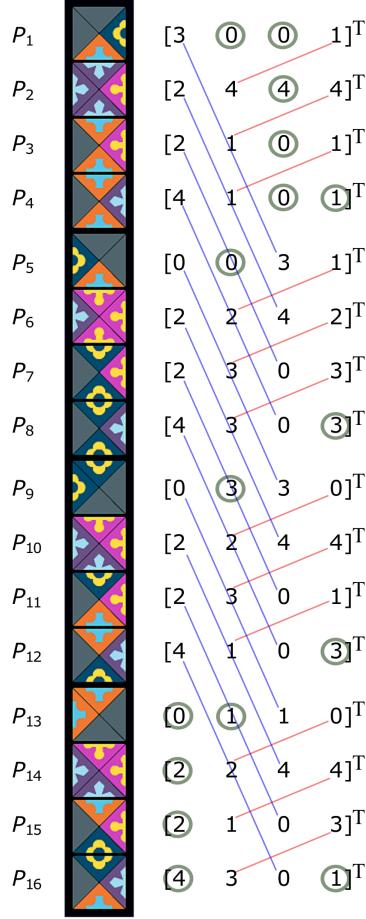


Figure 129: All pieces in their initial state on vectorized demo-game board. Line segments indicate differences  $\Delta$  (951),  $\bigcirc$  indicate edges on board boundary  $\beta$  (953). Entries are indices  $\ell$  identifying standard basis vectors  $e_\ell \in \mathbb{R}^L$  from Figure 128.

by a constant cardinality-2 vector  $\Delta_i$ , whose entries belong to  $\{-1, 0, 1\}$ , from a set  $\{\Delta_i \in \mathbb{R}^{4M}, i=1 \dots 2\sqrt{M}(\sqrt{M}-1)\}$ . Defining sparse constant wide matrix

$$\Delta \triangleq \begin{bmatrix} \Delta_1^T \\ \vdots \\ \Delta_{2\sqrt{M}(\sqrt{M}-1)}^T \end{bmatrix} \in \mathbb{R}^{2\sqrt{M}(\sqrt{M}-1) \times 4M} \quad (949)$$

then the desired constraint is

$$\Delta(\Xi \otimes I_4)\Pi P = \mathbf{0} \in \mathbb{R}^{2\sqrt{M}(\sqrt{M}-1) \times L} \quad (950)$$

For the demonstration game, the first twelve entries of  $\Delta$  correspond to blue line segments (leftmost) in Figure 129 while the twelve remaining entries correspond to red lines: for  $e_i \in \mathbb{R}^{64}$

$$\Delta = \begin{bmatrix} e_1^T - e_{19}^T \\ e_5^T - e_{23}^T \\ e_9^T - e_{27}^T \\ e_{13}^T - e_{31}^T \\ e_{17}^T - e_{35}^T \\ e_{21}^T - e_{39}^T \\ e_{25}^T - e_{43}^T \\ e_{29}^T - e_{47}^T \\ e_{33}^T - e_{51}^T \\ e_{37}^T - e_{55}^T \\ e_{41}^T - e_{59}^T \\ e_{45}^T - e_{63}^T \\ e_4^T - e_6^T \\ e_8^T - e_{10}^T \\ e_{12}^T - e_{14}^T \\ e_{20}^T - e_{22}^T \\ e_{24}^T - e_{26}^T \\ e_{28}^T - e_{30}^T \\ e_{36}^T - e_{38}^T \\ e_{40}^T - e_{42}^T \\ e_{44}^T - e_{46}^T \\ e_{52}^T - e_{54}^T \\ e_{56}^T - e_{58}^T \\ e_{60}^T - e_{62}^T \end{bmatrix} \in \mathbb{R}^{24 \times 64} \quad (951)$$

**game board boundary  $\beta$** 

Boundary of the square board must be colored grey. This means there are  $4\sqrt{M}$  boundary locations in  $\mathbb{R}^{4M \times L}$  (946) that must have value  $\mathbf{0}^T$ . Because  $(\Xi \otimes I_4)\Pi P \geq \mathbf{0}$ , these may all be lumped into one equality constraint

$$\beta^T (\Xi \otimes I_4)\Pi P \mathbf{1} = 0 \quad (952)$$

where  $\beta \in \mathbb{R}_+^{4M}$  is a sparse vector constant having entries in  $\{0, 1\}$  complementary to the known  $4\sqrt{M}$  zeros. For the demonstration game board Figure 129, for example,

$$\beta = [01100010001000110100000000000010100000000000011100100010001001]^T \in \mathbb{R}^{64} \quad (953)$$

**consolidating permutations  $\Phi$** 

By defining

$$\Phi \triangleq (\Xi \otimes I_4)\Pi \in \mathbb{R}^{4M \times 4M} \quad (954)$$

this square matrix becomes a structured permutation matrix replacing the product of permutation matrices. Then puzzle piece edge adjacency constraint (950) becomes

$$\Delta \Phi P = \mathbf{0} \in \mathbb{R}^{2\sqrt{M}(\sqrt{M}-1) \times L} \quad (955)$$

while game board boundary constraint (952) becomes

$$\beta^T \Phi P \mathbf{1} = 0 \in \mathbb{R} \quad (956)$$

Now partition composite permutation matrix variable  $\Phi$  into  $4 \times 4$  blocks

$$\Phi \triangleq \begin{bmatrix} \phi_{11} & \cdots & \phi_{1M} \\ \vdots & \ddots & \vdots \\ \phi_{M1} & \cdots & \phi_{MM} \end{bmatrix} \in \mathbb{R}^{4M \times 4M} \quad (957)$$

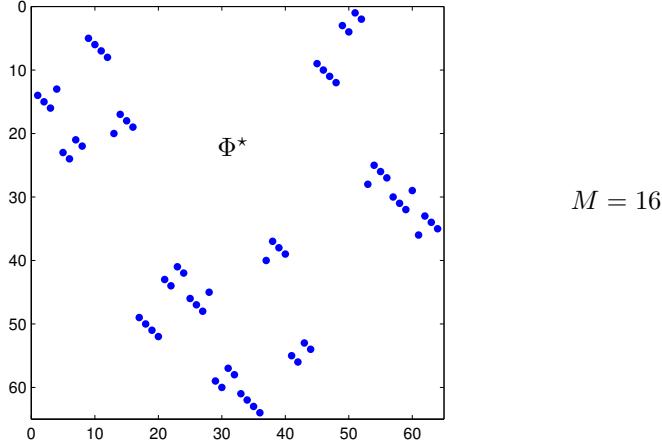


Figure 130: Sparsity pattern for composite permutation matrix  $\Phi^* \in \mathbb{R}^{4M \times 4M}$  representing solution from Figure 125b. Each four-point cluster represents a circulant permutation matrix from (947). Any  $M=16$ -piece solution may be verified on the [TOMY website](#).

where  $\Phi_{ij}^* \in \{0, 1\}$  because (947)

$$\phi_{ij}^* \in \{\mathbf{0}, \pi_1, \pi_2, \pi_3, \pi_4\} \subset \mathbb{R}^{4 \times 4} \quad (958)$$

An optimal composite permutation matrix  $\Phi^*$  is represented pictorially in Figure 130. Now we ask what are necessary conditions on  $\Phi^*$  at optimality:

- $4M$ -sparse<sup>4.78</sup> (cardinality-1 per row or column) and nonnegativity.
- Each column has one 1. Each row has one 1.
- Entries along each and every diagonal of each and every  $4 \times 4$  block  $\phi_{ij}^*$  are equal.
- Corner pair of  $2 \times 2$  submatrices on antidiagonal of each and every  $4 \times 4$  block  $\phi_{ij}^*$  are equal.

We want an objective function whose global optimum, when attained, certifies that the puzzle has been solved. Then, in terms of this  $\Phi$  partitioning (957), the Eternity II problem is a minimization of cardinality with optimal objective value  $8M$ .<sup>4.79</sup>

$$\begin{aligned}
 & \underset{\Phi \in \mathbb{R}^{4M \times 4M}}{\text{minimize}} \quad \sum_{i=1}^{4M} \|\Phi(i, :)^\top\|_0 + \|\Phi(:, i)\|_0 \\
 & \text{subject to} \quad \Delta \Phi P = \mathbf{0} \\
 & \quad \beta^\top \Phi P \mathbf{1} = 0 \\
 & \quad \Phi \mathbf{1} = \mathbf{1} \\
 & \quad \Phi^\top \mathbf{1} = \mathbf{1} \\
 & \quad (I_M \otimes R_d) \Phi (I_M \otimes R_d^\top) = (I_M \otimes S_d) \Phi (I_M \otimes S_d^\top) \\
 & \quad (I_M \otimes R_\phi) \Phi (I_M \otimes S_\phi^\top) = (I_M \otimes S_\phi) \Phi (I_M \otimes R_\phi^\top) \\
 & \quad (e_{121} \otimes I_4)^\top \Phi (e_{139} \otimes I_4) = \pi_3 \\
 & \quad \Phi \geq \mathbf{0}
 \end{aligned} \quad (959)$$

<sup>4.78</sup>Define *sparsity* as ratio of number of nonzero entries to matrix-dimension product. For matrices, the average number of nonzeros per row or column is easier to understand and likely to be small for typical LP problems, independent of the dimensions. —Michael Saunders

<sup>4.79</sup>A nonobvious method to transform cardinality minimization, in permutation problems, to rank minimization is disclosed in Example 4.7.0.0.3 with reference to Figure 118.

which is convex in the constraints where  $e_{121}, e_{139} \in \mathbb{R}^M$  are members of the standard basis representing mandatory piece  $P_{139}$  placement in the full game,<sup>4.80</sup> where

$$R_d \triangleq \begin{bmatrix} 1 & 0 & \mathbf{0} \\ & 1 & 0 \\ \mathbf{0} & & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 4}, \quad S_d \triangleq \begin{bmatrix} 0 & 1 & \mathbf{0} \\ & 0 & 1 \\ \mathbf{0} & & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 4} \quad (960)$$

$$R_\phi \triangleq \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \in \mathbb{R}^{2 \times 4}, \quad S_\phi \triangleq \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{2 \times 4} \quad (961)$$

and where  $\Phi \geq \mathbf{0}$  denotes entrywise nonnegativity. These matrices  $R$  and  $S$  enforce circulance.<sup>4.81</sup> Full game mandatory-piece rotation requires equality constraint  $\pi_3$ .

### permutation polyhedron

Constraints  $\Phi \mathbf{1} = \mathbf{1}$  and  $\Phi^T \mathbf{1} = \mathbf{1}$  and  $\Phi \geq \mathbf{0}$  confine  $\Phi$  to a permutation polyhedron (102) in  $\mathbb{R}^{4M \times 4M}$ ; which is, the convex hull of permutation matrices. The objective enforces minimal cardinality per row and column. Slicing the permutation polyhedron, by looking at a particular row or column subspace of  $\Phi$ , looks like intersection of a 1-norm ball with a nonnegative orthant. Cardinality 1 vectors reside at vertices of a one norm ball. (Figure 74)<sup>4.82</sup> Hence, the optimal objective is a sum of cardinalities 1.

Any vertex, of the permutation polyhedron, is a permutation matrix having minimal cardinality  $4M$ .<sup>4.83</sup> The feasible set of problem (959) is an intersection of the polyhedron with a number of hyperplanes. Feasible solutions exist that are not permutation matrices. But the intersection must contain a vertex of the permutation polyhedron because a solution  $\Phi^*$  cannot otherwise be a permutation matrix; such a solution is presumed to exist, so it must also be a vertex (extreme point)<sup>4.84</sup> of the intersection.

In vectorized variable, by §A.1.1 no.33, problem (959) is equivalent to

$$\begin{aligned} & \underset{\Phi \in \mathbb{R}^{4M \times 4M}}{\text{minimize}} \quad \sum_{i=1}^{4M} \|\Phi(i, :)^T\|_0 + \|\Phi(:, i)\|_0 \\ & \text{subject to} \quad (P^T \otimes \Delta) \text{vec } \Phi = \mathbf{0} \\ & \quad (P\mathbf{1} \otimes \beta)^T \text{vec } \Phi = 0 \\ & \quad (\mathbf{1}_{4M}^T \otimes I_{4M}) \text{vec } \Phi = \mathbf{1} \\ & \quad (I_{4M} \otimes \mathbf{1}_{4M}^T) \text{vec } \Phi = \mathbf{1} \\ & \quad (I_M \otimes R_d \otimes I_M \otimes R_d - I_M \otimes S_d \otimes I_M \otimes S_d) \text{vec } \Phi = \mathbf{0} \\ & \quad (I_M \otimes S_\phi \otimes I_M \otimes R_\phi - I_M \otimes R_\phi \otimes I_M \otimes S_\phi) \text{vec } \Phi = \mathbf{0} \\ & \quad (e_{139} \otimes I_4 \otimes e_{121} \otimes I_4)^T \text{vec } \Phi = \text{vec } \pi_3 \\ & \quad \text{vec } \Phi \succeq 0 \end{aligned} \quad (962)$$

whose optimal objective value is  $8M$ ; cardinality of permutation matrix  $\Phi^*$  is  $4M$ . With respect to an orthant,  $\succeq$  connotes entrywise nonnegativity (p.634). This problem is abbreviated:

$$\begin{aligned} & \underset{\Phi \in \mathbb{R}^{4M \times 4M}}{\text{minimize}} \quad \sum_{i=1}^{4M} \|\Phi(i, :)^T\|_0 + \|\Phi(:, i)\|_0 \\ & \text{subject to} \quad E \text{vec } \Phi = \tau \\ & \quad \text{vec } \Phi \succeq 0 \end{aligned} \quad (963)$$

<sup>4.80</sup> meaning that piece  $P_{139}$  (numbered 139 by Monckton) must be placed in cell 121 on the board (Figure 126) with rotation  $\pi_3$  (p.311).

<sup>4.81</sup>Since  $\mathbf{0}$  is the trivial circulant matrix, application is democratic over all blocks  $\phi_{ij}$ .

<sup>4.82</sup>This means: each vertex of the permutation polyhedron, in isometrically isomorphic  $\mathbb{R}^{16M^2}$ , is coincident with a vertex of  $8M$  1-norm balls in  $4M$ -dimensional subspaces.

<sup>4.83</sup>but maximal Frobenius norm (p.322).

<sup>4.84</sup>Vertex means zero-dimensional exposed face (§2.6.1.0.1); intersection with a strictly supporting hyperplane. There can be no further intersection with a feasible affine subset that would enlarge that face; *id est*, a vertex of the permutation polyhedron persists in the feasible set.

where  $E \in \mathbb{R}^{17+2L\sqrt{M}(\sqrt{M}-1)+8M+13M^2 \times 16M^2}$  is highly sparse having 4,784,144 nonzero entries in  $\{-1, 0, 1\}$ .

$$\dim E = 864,593 \times 1,048,576 \quad (964)$$

- Any feasible binary solution is minimal cardinality and *vice versa* because it is a vertex of the feasible set. (§2.3.2.0.4)

But number of equality constraints is too large for contemporary binary solvers.<sup>4.85</sup> So again, we reformulate the problem:

### canonical Eternity II

Because each block  $\phi_{ij}$  of  $\Phi$  (957) is optimally circulant, comprising four permutation matrices (958) uniquely identifiable by their first column (947), we may take as variable every fourth column of  $\Phi$ :

$$\tilde{\Phi} \triangleq [\Phi(:, 1) \ \Phi(:, 5) \ \Phi(:, 9) \ \dots \ \Phi(:, 4M-3)] \in \mathbb{R}^{4M \times M} \quad (965)$$

where  $\tilde{\Phi}_{ij} \in \{0, 1\}$ . Then, for  $e_i \in \mathbb{R}^4$

$$\Phi = (\tilde{\Phi} \otimes e_1^T) + (I_M \otimes \pi_4)(\tilde{\Phi} \otimes e_2^T) + (I_M \otimes \pi_3)(\tilde{\Phi} \otimes e_3^T) + (I_M \otimes \pi_2)(\tilde{\Phi} \otimes e_4^T) \in \mathbb{R}^{4M \times 4M} \quad (966)$$

This formula describes replication (+), columnar upsampling & shifting ( $e_i \in \mathbb{R}^4$ ), and rotation ( $\pi_i \in \mathbb{R}^{4 \times 4}$ ) of  $\tilde{\Phi}$ . By §A.1.1 no.45 and no.46

$$\begin{aligned} \text{vec } \Phi &= (I_M \otimes e_1 \otimes I_{4M} + I_M \otimes e_2 \otimes I_M \otimes \pi_4 + I_M \otimes e_3 \otimes I_M \otimes \pi_3 + I_M \otimes e_4 \otimes I_M \otimes \pi_2) \text{vec } \tilde{\Phi} \\ &\triangleq Y \text{vec } \tilde{\Phi} \in \mathbb{R}^{16M^2} \end{aligned} \quad (967)$$

where  $Y \in \mathbb{R}^{16M^2 \times 4M^2}$ . Because three out of every four rows (per consecutive quadruple adjacent rows of  $\tilde{\Phi}$ ) equal  $\mathbf{0}^T$ , permutation polyhedron (102) demands that each quadruple and each column sum to 1: respectively,  $(I_M \otimes \mathbf{1}_4^T)\tilde{\Phi}\mathbf{1} = \mathbf{1}$  and  $\tilde{\Phi}^T\mathbf{1} = \mathbf{1}$  where  $\tilde{\Phi}$  is now variable and optimally binary. By substitution of columnar subsampled matrix  $\tilde{\Phi}$  (965) for permutation matrix  $\Phi$  (954), circulance constraints in  $R$  and  $S$  (which are most numerous) may be dropped from Eternity II problem (959) because circulance of  $\phi_{ij}$  is built into  $\Phi$ -reconstruction formula (966). We are left with a feasibility problem equivalent to (959), for  $e_{121}, e_{139} \in \mathbb{R}^M$

$$\begin{aligned} &\text{find } \tilde{\Phi} \in \mathbb{B}^{4M \times M} \\ &\text{subject to } \Delta \Phi P = \mathbf{0} \\ &\quad \beta^T \Phi P \mathbf{1} = 0 \\ &\quad (I_M \otimes \mathbf{1}_4^T)\tilde{\Phi}\mathbf{1} = \mathbf{1} \\ &\quad \tilde{\Phi}^T\mathbf{1} = \mathbf{1} \\ &\quad (e_{121} \otimes I_4)^T \Phi (e_{139} \otimes I_4) = \pi_3 \end{aligned} \quad (968)$$

where  $\Delta \in \mathbb{R}^{2\sqrt{M}(\sqrt{M}-1) \times 4M}$  (949) (identifying adjacent edges) is evaluated in (951), initial piece placement  $P \in \mathbb{R}^{4M \times L}$  is defined in (944) and enumerated in Figure 128,  $\beta \in \mathbb{R}_+^{4M}$  defining a game board boundary in Figure 129 has corresponding value (953), and where  $\pi_3$  (947) determines mandatory-piece rotation. Thus, Eternity II (962) is equivalently

---

<sup>4.85</sup> Saunders' program `lusol` can reduce that number to 797,508 constraints by eliminating linearly dependent rows of matrix  $E$ , but that is not enough to overcome numerical issues with the best solvers.

transformed

$$\begin{aligned}
 & \underset{\tilde{\Phi} \in \mathbb{R}^{4M \times M}}{\text{minimize}} && \sum_{i=1}^{4M} \|\tilde{\Phi}(i, :)^T\|_0 + \sum_{j=1}^M \|\tilde{\Phi}(:, j)\|_0 \\
 & \text{subject to} && (P^T \otimes \Delta)Y \text{ vec } \tilde{\Phi} = \mathbf{0} \\
 & && (P\mathbf{1} \otimes \beta)^T Y \text{ vec } \tilde{\Phi} = 0 \\
 & && (\mathbf{1}_M^T \otimes I_M \otimes \mathbf{1}_4^T) \text{ vec } \tilde{\Phi} = \mathbf{1} \\
 & && (I_M \otimes \mathbf{1}_{4M}^T) \text{ vec } \tilde{\Phi} = \mathbf{1} \\
 & && (e_{139} \otimes e_{121} \otimes I_4)^T \text{ vec } \tilde{\Phi} = \pi_3 e_1 \\
 & && \text{vec } \tilde{\Phi} \succeq 0
 \end{aligned} \tag{969}$$

whose optimal objective value is  $2M$  since optimal cardinality of  $\tilde{\Phi}^*$  (with entries in  $\{0, 1\}$ ) is  $M$ , where matrix constant  $Y$  maps subsampled  $\tilde{\Phi}$  to  $\Phi$  via (967), and where  $e_1 \in \mathbb{R}^4$ . In abbreviation of reformulation (969)

$$\begin{aligned}
 & \underset{\tilde{\Phi} \in \mathbb{R}^{4M \times M}}{\text{minimize}} && \sum_{i=1}^{4M} \|\tilde{\Phi}(i, :)^T\|_0 + \sum_{j=1}^M \|\tilde{\Phi}(:, j)\|_0 \\
 & \text{subject to} && \tilde{E} \text{ vec } \tilde{\Phi} = \tilde{\tau} \\
 & && \text{vec } \tilde{\Phi} \succeq 0
 \end{aligned} \tag{970}$$

number of equality constraints is now 11,077; an order of magnitude fewer constraints than (963). Sparse  $\tilde{E} \in \mathbb{R}^{5+2L\sqrt{M}(\sqrt{M}-1)+2M \times 4M^2}$  replaces matrix  $E$ . Number of columns has also been reduced, down from more than a million:

$$\dim \tilde{E} = 11,077 \times 262,144 \tag{971}$$

But this dimension remains out of reach of most highly regarded academic and commercial binary solvers; especially disappointing insofar as sparsity of  $\tilde{E}$  is high with 1,503,732 nonzero entries in  $\{-1, 0, 1, 2\}$ ; element {2} arising only in the  $\beta$  constraint which is soon to disappear after presolving.

### presolving: game board's edge

Any process of discarding rows and columns, prior to numerical optimization, is called *presolving*. The constraint in  $\beta$ , which zeroes the board at its edges, has all positive coefficients. The zero sum means that all  $\text{vec } \tilde{\Phi}$  entries, corresponding to nonzero entries in row vector  $(P\mathbf{1} \otimes \beta)^T Y$ , must be zero. For the full game, this means we may immediately eliminate 57,840 variables from 262,144. After zero-row and dependent-row (two) removal,

$$\dim \tilde{E} \rightarrow 10,054 \times 204,304 \tag{972}$$

with entries in  $\{-1, 0, 1\}$ .

### geometric presolver: polyhedral cone theory

Eternity II problem (970) constraints are interpretable in the language of convex cones: The columns of matrix  $\tilde{E}$  constitute a set of generators for a pointed (§2.12.2.2) polyhedral cone

$$\mathcal{K} = \{\tilde{E} \text{ vec } \tilde{\Phi} \mid \text{vec } \tilde{\Phi} \succeq 0\} \tag{973}$$

Even more intriguing is the observation: vector  $\tilde{\tau}$  resides on that polyhedral cone's boundary.<sup>4.86</sup> (§2.13.4.2.4) We may apply techniques from §2.13.5 to prune generators not belonging to the smallest face of that cone, to which  $\tilde{\tau}$  belongs, because generators of

---

<sup>4.86</sup>This observation applies equally well to cones generated by both (971) and (972). And  $\tau$  is on the boundary of the polyhedral cone generated by  $E$  (964).

that smallest face must hold a minimal cardinality solution. Matrix dimension is thereby reduced:<sup>4.87</sup>

Designate  $\mathcal{I}$  as the set of all surviving column indices of  $\tilde{E}$  from  $4M^2 = 262,144$  columns:

$$\mathcal{I} \subset \{1 \dots 4M^2\} \quad (974)$$

The  $i^{\text{th}}$  column  $\tilde{E}(:, i)$  of matrix  $\tilde{E}$  belongs to the smallest face  $\mathcal{F}$  of  $\mathcal{K}$  that contains  $\tilde{\tau}$  if and only if

$$\begin{aligned} & \text{find} && \text{vec } \tilde{\Phi} \in \mathbb{R}^{\dim \mathcal{I}}, \mu \in \mathbb{R} \\ & \text{subject to} && \mu \tilde{\tau} - \tilde{E}(:, i) = \tilde{E} \text{ vec } \tilde{\Phi} \\ & && \text{vec } \tilde{\Phi} \succeq 0 \end{aligned} \quad (380)$$

is feasible. By a transformation of [Saunders](#), this linear feasibility problem is the same as

$$\begin{aligned} & \text{find} && \text{vec } \tilde{\Phi} \in \mathbb{R}^{\dim \mathcal{I}}, \mu \in \mathbb{R} \\ & \text{subject to} && \tilde{E} \text{ vec } \tilde{\Phi} = \mu \tilde{\tau} \\ & && \text{vec } \tilde{\Phi} \succeq 0 \\ & && (\text{vec } \tilde{\Phi})_i \geq 1 \end{aligned} \quad (975)$$

A minimal cardinality solution to Eternity II (970) implicitly constrains  $\tilde{\Phi}^*$  to be binary. So this test (975) of membership to  $\mathcal{F}(\mathcal{K} \ni \tilde{\tau})$  may be tightened to a test of  $(\text{vec } \tilde{\Phi})_i = 1$ ; *id est*, for  $i = 1 \dots \dim \mathcal{I} = 1 \dots 204,304$  distinct linear feasibility problems

$$\begin{aligned} & \text{find} && \text{vec } \tilde{\Phi} \in \mathbb{R}^{\dim \mathcal{I}} \\ & \text{subject to} && \tilde{E} \text{ vec } \tilde{\Phi} = \tilde{\tau} \\ & && \text{vec } \tilde{\Phi} \succeq 0 \\ & && (\text{vec } \tilde{\Phi})_i = 1 \end{aligned} \quad (976)$$

whose feasible set is a proper subset of that in (975). Real variable  $\mu$  can be set to 1 because if it must not be, then feasible  $(\text{vec } \tilde{\Phi})_i = 1$  could not be feasible to Eternity II (970).

If infeasible here in (976), then the only choice remaining for  $(\text{vec } \tilde{\Phi})_i$  is 0; meaning, column  $\tilde{E}(:, i)$  may be discarded but only after all columns have been tested. This tightened problem (976) therefore tells us two things when feasible:  $\tilde{E}(:, i)$  belongs to the smallest face of  $\mathcal{K}$  that contains  $\tilde{\tau}$ , and  $(\text{vec } \tilde{\Phi})_i$  constitutes a nonzero vertex-coordinate of permutation polyhedron (102). After presolving via this conic pruning method (with subsequent zero-row and dependent-row removal),

$$\dim \tilde{E} \rightarrow 7,362 \times 150,638 \quad (977)$$

Entries in  $\text{vec } \tilde{\Phi}$ , corresponding to discarded columns of  $\tilde{E}$ , are optimally 0. But now  $\tilde{\tau}$  resides relatively interior to the polyhedral cone (973) generated by this reduction  $\tilde{E}$ . Its binary nature is evident in Hu's depiction of our reduced "A" matrix in Figure 131.

**c.i. presolver Eternity II: Generators of smallest face are conically independent**  
 Matrix  $\tilde{E}$  now accounts for the board's edge and holds what remains after discard of all generators not in the smallest face  $\mathcal{F}$  of cone  $\mathcal{K}$  that contains  $\tilde{\tau}$ . To further prune

---

<sup>4.87</sup>Column elimination can be quite dramatic but is dependent upon problem geometry. By method of convex cones, we will discard 53,666 more columns via Saunders' [pdco](#); a total of 111,506 columns will have been removed from 262,144. Following dependent-row removal via [lusol](#), dimension of  $\tilde{E}$  becomes  $7,362 \times 150,638$ .

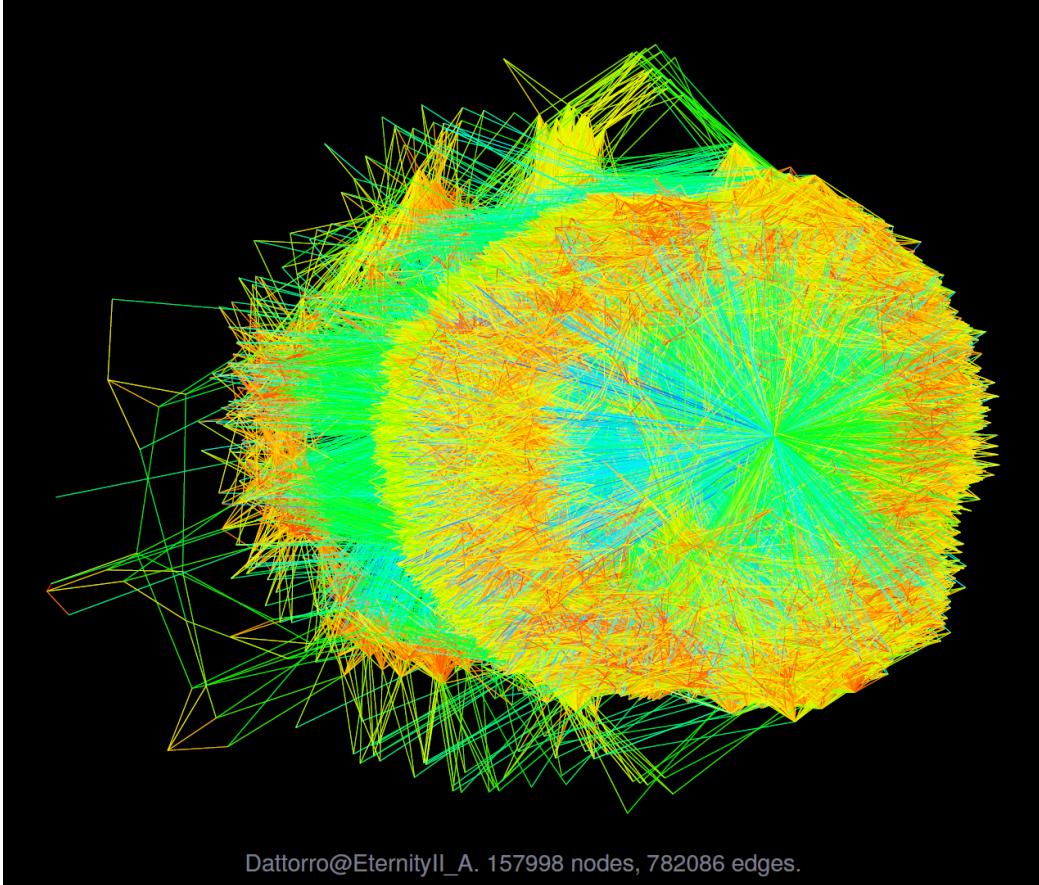


Figure 131: Directed graph of *adjacency matrix* for  $\tilde{E}$  (977) ( $\tilde{E} = A$  in [415]) representing reduced equality constraint in Eternity II problem. “Movie” in [231] shows realization in three dimensions; color corresponding to line-segment length. (Realization by Yifan Hu.)

generators relatively interior to that smallest face, we may subsequently test for conic dependence as described in §2.10: for  $i=1 \dots \dim \mathcal{I} = 1 \dots 150,638$

$$\begin{aligned} &\text{find } \text{vec } \tilde{\Phi} \in \mathbb{R}^{\dim \mathcal{I}} \\ &\text{subject to } \tilde{E} \text{ vec } \tilde{\Phi} = \tilde{E}(:, i) \\ &\quad \text{vec } \tilde{\Phi} \succeq 0 \\ &\quad (\text{vec } \tilde{\Phi})_i = 0 \end{aligned} \tag{283}$$

If feasible, then column  $\tilde{E}(:, i)$  is a conically dependent generator of the smallest face and must be discarded from matrix  $\tilde{E}$  before proceeding with test of remaining columns.

Generators interior to a smallest face could provide a lower cardinality solution, so it might be imprudent to prune. It turns out, for Eternity II: generators of the smallest face, previously found via (976), comprise a minimal set; *id est*, (283) is never feasible; no more columns of  $\tilde{E}$  can be discarded.<sup>4.88</sup>

$$m \times \dim \mathcal{I} \triangleq \dim \tilde{E} = 7,362 \times 150,638 \tag{977}$$

---

<sup>4.88</sup>One cannot help but notice a binary selection of variable by tests (976) and (283): Geometrical test (976) (smallest face) checks feasibility of vector entry 1 while geometrical test (283) (conic independence) checks feasibility of 0. Changing 1 to 0 in (976) is always feasible for Eternity II.

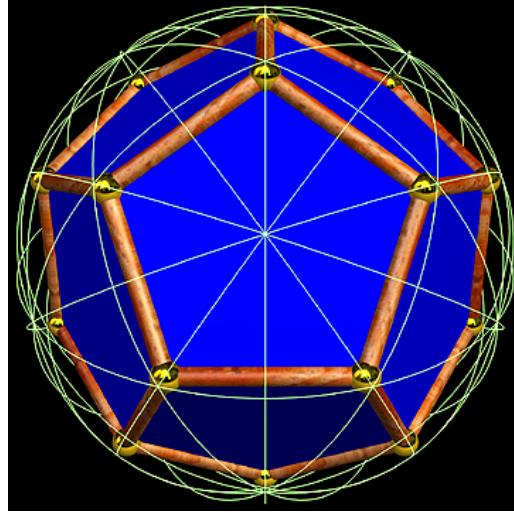


Figure 132: Polyhedron vertices  $\bullet$  inscribed on sphere skeleton in  $\mathbb{R}^3$  for visualization of permutation matrices in abstract isomorphic space. Vertices represent matrices, of the same dimension, all equidistant from origin. Sphere about origin represents level set at maximum of simple quadratic  $x^T x$  where vertices intersect sphere. Permutation matrices are represented by those vertices in nonnegative orthant. If sphere expands, intersection with polyhedron becomes empty. (Drawing by Robert Austin using [Stella4D](#).)

Successive reductions of  $E$  and  $\tau$  can be found on [Wikimization](#) [415] in MATLAB format.  $\square$

Incorporating more Clue Pieces, provided by Monckton, makes the Eternity II problem harder in the sense that solution set is diminished; the target gets smaller.<sup>4.89</sup>

#### 4.7.0.0.16 Example. Eternity II - affinity for maximization.

Reversing tack on Example 4.7.0.0.15, Eternity II optimization resembles Figure 33a (not (b)) because variable  $\tilde{\Phi}$  is implicitly bounded above by design;  $\mathbf{1} \succeq \text{vec } \tilde{\Phi}$  by confinement of  $\Phi$  in (959) to the permutation polyhedron (102), for  $i=1 \dots \dim \mathcal{I}=1 \dots 150,638$

$$\begin{aligned} 1 = & \underset{\tilde{\Phi}}{\text{maximize}} \quad (\text{vec } \tilde{\Phi})_i \\ & \text{subject to} \quad \tilde{E} \text{ vec } \tilde{\Phi} = \tilde{\tau} \\ & \text{vec } \tilde{\Phi} \succeq 0 \end{aligned} \tag{978}$$

Unity is always attainable, by (976). By (965) this means (§4.6.1.4)

$$\begin{aligned} M = & \underset{\substack{y(\tilde{\Phi}), \tilde{\Phi}}}{\text{maximize}} \quad (\mathbf{1} - y)^T \text{vec } \tilde{\Phi} \quad \underset{\tilde{\Phi}}{\text{maximize}} \quad \|\text{vec } \tilde{\Phi}\|_{\dim \mathcal{I}} \\ & \text{subject to} \quad \tilde{E} \text{ vec } \tilde{\Phi} = \tilde{\tau} \quad \equiv \quad \text{subject to} \quad \tilde{E} \text{ vec } \tilde{\Phi} = \tilde{\tau} \\ & \text{vec } \tilde{\Phi} \succeq 0 \quad \text{vec } \tilde{\Phi} \succeq 0 \end{aligned} \tag{979}$$

where

$$y = \mathbf{1} - \nabla \|\text{vec } \tilde{\Phi}\|_{\dim \mathcal{I}} \tag{842}$$

---

<sup>4.89</sup>But given the four clues provided, our geometric presolver (p.318) produces a 15% smaller face; a total very nearly half the 262,144 columns can be proven to correspond to 0 coefficients.

is a direction vector from the cardinality minimization technique of convex iteration in §4.6.1.1 and where  $\|\text{vec } \tilde{\Phi}\|_{\dim_M^{\mathcal{I}}}$  is a  $k$ -largest norm (§3.2.2.1,  $k = M$ ). When upper bound  $M$  in (979) is met, solution  $\text{vec } \tilde{\Phi}^*$  will be optimal for Eternity II because it must then be a Boolean vector with minimal cardinality  $M$ .

Maximization of convex function  $\|\text{vec } \tilde{\Phi}\|_{\dim_M^{\mathcal{I}}}$  (monotonic on  $\mathbb{R}_+^{\dim \mathcal{I}}$ ) is not a convex problem, though the constraints are convex. [343, §32] Geometrical visualization of this problem formulation is clear. We therefore choose to work with a complementary direction vector  $z$ , in what follows, in predilection for a mental picture of convex function maximization.

### complementary direction vector is optimal solution of Eternity II

Instead of solving (979), which is difficult, we propose iterating a convex problem sequence: for  $\mathbf{1} - y \leftarrow z$

$$\begin{aligned} & \underset{\text{vec } \tilde{\Phi} \in \mathbb{R}^{\dim \mathcal{I}}}{\text{maximize}} && z^T \text{vec } \tilde{\Phi} \\ & \text{subject to} && \tilde{E} \text{vec } \tilde{\Phi} = \tilde{\tau} \\ & && \text{vec } \tilde{\Phi} \succeq 0 \end{aligned} \quad (980)$$

$$\begin{aligned} & \underset{z \in \mathbb{R}^{\dim \mathcal{I}}}{\text{maximize}} && z^T \text{vec } \tilde{\Phi}^* \\ & \text{subject to} && 0 \preceq z \preceq \mathbf{1} \\ & && z^T \mathbf{1} = M \end{aligned} \quad (531)$$

Variable  $\tilde{\Phi}$  is implicitly bounded above at  $\mathbf{1}$  by design of  $\tilde{E}$ . A globally optimal complementary direction vector  $z^*$  will always exactly match an optimal solution  $\text{vec } \tilde{\Phi}^*$  for convex iteration of any problem formulated as maximization of a Boolean variable: here we have

$$z^{*T} \text{vec } \tilde{\Phi}^* \triangleq M \quad (981)$$

Because  $z^* = \text{vec } \tilde{\Phi}^*$ , Eternity II can be equivalently formulated as maximization of a convex quadratic instead:

$$\begin{aligned} & \underset{\text{vec } \tilde{\Phi} \in \mathbb{R}^{\dim \mathcal{I}}}{\text{maximize}} && (\text{vec } \tilde{\Phi})^T \text{vec } \tilde{\Phi} \\ & \text{subject to} && \tilde{E} \text{vec } \tilde{\Phi} = \tilde{\tau} \\ & && \text{vec } \tilde{\Phi} \succeq 0 \end{aligned} \quad (982)$$

a nonconvex problem but requiring no convex iteration. The optimal objective is known:  $(\text{vec } \tilde{\Phi}^*)^T \text{vec } \tilde{\Phi}^* = \|\tilde{\Phi}^*\|_F^2 = \mathbf{1}^T \text{vec } \tilde{\Phi}^* = M$  with  $\text{vec } \tilde{\Phi}^*$  binary and cardinality- $M$  attained at a vertex of the permutation polyhedron (p.316). (Figure 132)

### rumination

If it were possible to form a nullspace basis  $Z$  for  $\tilde{E}$ , of about equal sparsity such that

$$\text{vec } \tilde{\Phi} = Z\xi + \text{vec } \tilde{\Phi}_p \quad (117)$$

then a problem formulation equivalent to (982)

$$\begin{aligned} & \underset{\xi}{\text{maximize}} && (Z\xi + \text{vec } \tilde{\Phi}_p)^T (Z\xi + \text{vec } \tilde{\Phi}_p) \\ & \text{subject to} && Z\xi + \text{vec } \tilde{\Phi}_p \succeq 0 \end{aligned} \quad (983)$$

might invoke optimality conditions as obtained in [223, thm.8].  $\square$

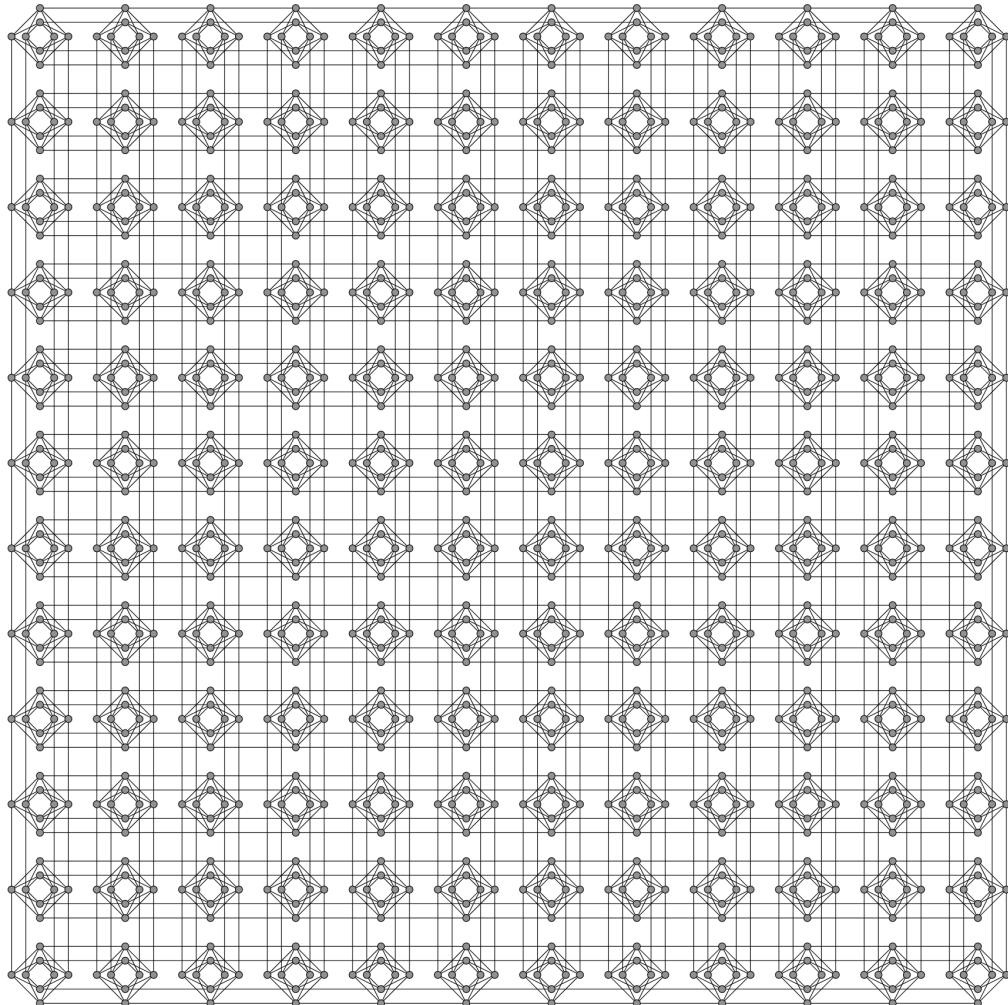


Figure 133:  $N \times N = 12 \times 12$  *Chimera* topology for D:Wave 1152-qubit chip architecture illustrating 3,360 ( $= 16 \times 12 \times 12 + 2 \times 11 \times 4 \times 12$ ) couplers by line segments between qubits (lines cross without intersection). Coupled qubits are *neighbors*, but distance is not preserved by this map. (Drawing by Diane Carr.)

## 4.8 Quantum optimization

*There was a time when the newspapers said that only twelve men understood the theory of relativity. I don't believe there ever was such a time. . . . a lot of people kind of understood the theory of relativity in some way or other, but more than twelve. On the other hand, I think I can safely say that nobody understands quantum mechanics.*

—Richard Feynman, 1964

A superconducting *quantum annealer* is the physical embodiment of an optimizer that globally minimizes a hypersurface whose modes increase factorially with dimension.

*Note that this architecture is very different from conventional computing. The processor has no large areas of memory (cache), rather each qubit has a tiny piece of memory of its own. In fact, the chip is architected more like a biological brain than the common ‘Von Neumann’ architecture of a conventional silicon processor. One can think of the qubits as being like neurons, and the couplers as being like synapses that control the flow of information between those neurons.*

—[dwavesys.com, §1.3]

A quantum annealer is unlike a von Neumann computer architecture insofar as it does not solve equations, there are no conditionally executable instructions, one *qubit* (the quantum analogue to bit) can be in the two binary states at once,<sup>4.90</sup> and qubit values may not be set by a programmer [101, §2]. There is no clock in a quantum annealer which operates at a temperature colder than outer space: near 0° Kelvin. The first commercially available quantum annealer was delivered in 2011.<sup>4.91</sup> Even though its magnetic superconducting *niobium* qubits are etched on a *silicon substrate*, a chip, the **D:Wave** quantum annealer is actually the first analog computer of its kind.<sup>4.92</sup>

*Ising’s spin model* [46, §2.1] [348, p.297] is a measure of molecular energy for a magnetic material, for bipolar binary  $s \in \mathbb{B}_{\pm}^n = \{-1, 1\}^n$

$$E(s) = \frac{1}{2} \langle J, ss^T \rangle + \langle h, s \rangle \quad (984)$$

Given *applied field strength*  $h \in \mathbb{R}^n$  and *interaction field strength*  $J \in \mathbb{S}_h^n$ , a quantum annealer minimizes this energy  $E$  which is always bounded because vector variable  $s$  is bounded above and below.

A graph of the D:Wave  $N \times N$  *Chimera* topology ( $N=12$ ) is represented in Figure 133; a *neighboring qubit* topology. Hollow matrix  $J$  represents *coupling* that occurs among physically neighboring qubits. Scalar  $\frac{1}{2}$  accounts for bidirectional coupling implied by  $J$  matrix symmetry. Coupling, which is an application of *entanglement* in quantum physics, can be controlled only for physically neighboring qubits. Increasing number of neighbors is therefore of practical importance. [45] Effective coupling of distant qubits is implemented by replicating qubits redundantly. [101, §3.4] As rule of thumb, complete coupling of  $n$  qubits (highest density  $J$ ) would leave  $O(\sqrt{n})$  qubits available.<sup>4.93</sup>

<sup>4.90</sup> the qubit’s *superposition* state.

<sup>4.91</sup> It is not capable of solving Eternity II in 2016 because of qubits insufficient in number and *coupling*.

<sup>4.92</sup> At present, there are two emergent technologies for harnessing quantum phenomena: *adiabatic model* (analog annealer) and *gate model* (analogue to Boolean logic gates of digital computers).

<sup>4.93</sup> 1152-qubit architecture machines, having 3,360 physical couplers, became available in 2015 for \$10M USD. In 2016, 2048-qubit chips (6,016 couplers) were announced. If qubit growth continues following *Rose’s law*, we should see million-qubit chips in 2025. Given  $n$  qubits, complete coupling requires  $n(n-1)/2$  couplers. Insufficient coupler population, not qubits, will become the bottleneck. In the near term, innovating a three-dimensional qubit topology would accelerate ratio of coupler to qubit growth.

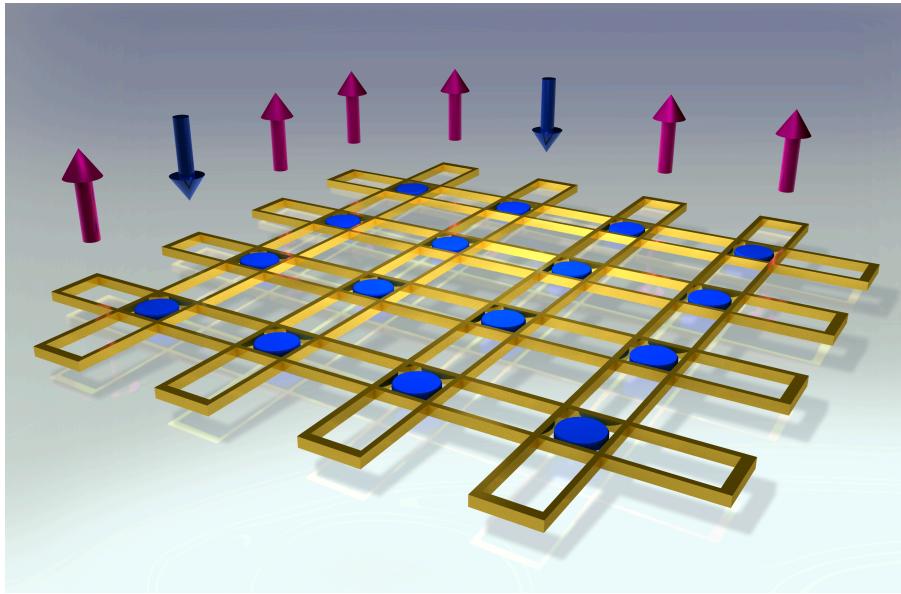


Figure 134: Chimera circuit chip layout abstract, topological dimension  $N=1$  illustrated. Eight qubits comprise hollow slabs  $\square$  whereas couplers are represented by sixteen blue discs  $\bullet$ . Chip layout is dual to graph topology in Figure 133. Up/down arrows connote final qubit states.

Chimera provides  $8N^2$  qubits having  $N(24N-8)$  physical couplers;  $\frac{\text{coupler}}{\text{qubit}} = 3 - \frac{1}{N}$  approaches three couplers per qubit on average, as topological dimension  $N$  increases, although it has no complete circuit of three qubits. Because couplers are bidirectional, each qubit effectively sees twice the coupler/qubit ratio:<sup>4.94</sup>

$$6 - \frac{2}{N} \text{ couplings/qubit} \quad (985)$$

For  $1 \leq N \leq 2$ , this coupling number is exact. For  $N > 2$ , this real number should be regarded as average number of couplings seen by a qubit. Physical layout of Chimera reveals a duality with graph topology: In chip layout Figure 134, physical qubits are represented by hollow slabs  $\square$  and couplers by discs  $\bullet$ . But in the topological graph in Figure 133, qubits are represented by discs  $\bullet$  and couplers by line segments  $—$ .

The D-Wave machine performs physical, not simulated, annealing. The system is initialized to a superposition (a  $2^n$  simultaneity) of all possible states [456, p.2] by application of a globally transverse magnetic field [204, slide 8/45]. At its outset, the energy hypersurface appears globally convex but settles into the Ising model after about  $20\mu\text{s}$  [433] with 2015 technology.<sup>4.95</sup> A globally optimal solution cannot be guaranteed because present understanding of the quantum annealing process is nondeterministic. To increase probability of finding a globally optimal solution, the same problem is sequentially executed thousands of times on the quantum annealer. The minimum, from each run, becomes a sample in proximity to the global minimum of binary quadratic function (984). (Sampling is necessary because successive minima can be offset by as much as a few percent from the global minimum.)

<sup>4.94</sup>Complete coupling is impossible with current technology; it would require an  $\frac{8N^2(8N^2-1)}{2}$  line-segment topology: were  $\frac{\text{coupler}}{\text{qubit}} = \frac{8N^2-1}{2}$ , each qubit would effectively see  $8N^2 - 1$  couplings.

<sup>4.95</sup>2011 saw 128-qubit machines with settling time at about  $75\mu\text{s}$  [246].

By change of variable, for binary  $q \in \mathbb{B}^n = \{0, 1\}^n$

$$s \leftarrow 2q - \mathbf{1} \quad (986)$$

the resulting *quadratic unconstrained binary optimization* (QUBO)

$$2 \underset{q \in \{0, 1\}^n}{\text{minimize}} \langle J, qq^T \rangle + \langle h - J\mathbf{1}, q \rangle \quad (987)$$

remains an equivalent energy minimization whose constant term  $\mathbf{1}^T(J\mathbf{1}\frac{1}{2} - h)$  is ignored. Whereas

$$\delta(ss^T) = \mathbf{1}, \quad \delta(qq^T) = q \quad (988)$$

this latter equality in  $q$  means that QUBO (987) is the same as

$$2 \underset{q \in \{0, 1\}^n}{\text{minimize}} \langle J, qq^T \rangle + \langle \delta(h - J\mathbf{1}), qq^T \rangle = 2 \underset{q \in \{0, 1\}^n}{\text{minimize}} \langle J + \delta(h - J\mathbf{1}), qq^T \rangle \quad (989)$$

Coefficient matrix  $J + \delta(h - J\mathbf{1})$  can be indefinite.

To abstract problem formulation away from the machine a little more (to simplify presentation), a QUBO shall be generalized

$$\underset{q \in \{0, 1\}^n}{\text{minimize}} q^T B q + a^T q \quad (990)$$

where  $a \in \mathbb{R}^n$  is a coefficient vector and where hollow matrix  $B \in \mathbb{R}^{n \times n}$  (comprising quadratic coefficients) is not necessarily symmetric; its main diagonal may be assumed  $\mathbf{0}$  and its lower triangular part empty.

#### 4.8.1 quantum gap maximization by linear programming

To further increase probability of finding a globally optimal solution, the discrete *gap* between optimal objective and least suboptimal objective is maximized by problem design [45] (992); by discriminating coefficients as in Example 4.8.2.0.1 and Example 4.8.2.0.2. D:Wave quantum annealer coefficient quantization is coarse, encoded by application of an external magnetic field whose resolution is about 4 or 5 bits over  $[-2, 2]$ .

Because the Eternity II puzzle (§4.7.0.0.15) can be formulated as a permutation problem, it is of interest to express a permutation polyhedron constraint (p.316). To illustrate realization of just one row of a permutation matrix in QUBO form (990), consider an  $n$ -qubit vector  $q$  that is allowed to have only one nonzero; *id est*, a discrete impulse ([a.k.a Kronecker delta](#)) over a vector of qubits. In other words, we need to translate this program

$$\begin{array}{ll} \text{find } & q \in \mathbb{B}^n \\ \text{subject to } & q^T \mathbf{1} = 1 \end{array} \quad (991)$$

into a QUBO. First we analyze a three-qubit case, then generalize to  $n$  in (994):

**quantum impulse:** . . .

$q_1$	$q_2$	$q_3$	$B_{12}q_1q_2 + B_{13}q_1q_3 + B_{23}q_2q_3 + a_1q_1 + a_2q_2 + a_3q_3$
desirable			
0	0	1	$a_3$
0	1	0	$a_2$
1	0	0	$a_1$
undesirable			
0	0	0	0
0	1	1	$B_{23} + a_2 + a_3$
1	0	1	$B_{13} + a_1 + a_3$
1	1	0	$B_{12} + a_1 + a_2$
1	1	1	$B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3$

Coefficients  $B$  and  $a$  (990) are selected by solution to a linear program whose undesirable objectives always exceed the objective for each and every desirable state:

$$\begin{aligned}
 & \underset{B, a, \text{gap}}{\text{maximize}} && \text{gap} \\
 & \text{subject to} && 0 \geq a_3 + \text{gap} \\
 & && 0 \geq a_2 + \text{gap} \\
 & && 0 \geq a_1 + \text{gap} \\
 & && B_{23} + a_2 + a_3 \geq a_3 + \text{gap} \\
 & && B_{23} + a_2 + a_3 \geq a_2 + \text{gap} \\
 & && B_{23} + a_2 + a_3 \geq a_1 + \text{gap} \\
 & && B_{13} + a_1 + a_3 \geq a_3 + \text{gap} \\
 & && B_{13} + a_1 + a_3 \geq a_2 + \text{gap} \\
 & && B_{13} + a_1 + a_3 \geq a_1 + \text{gap} \\
 & && B_{12} + a_1 + a_2 \geq a_3 + \text{gap} \\
 & && B_{12} + a_1 + a_2 \geq a_2 + \text{gap} \\
 & && B_{12} + a_1 + a_2 \geq a_1 + \text{gap} \\
 & && B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3 \geq a_3 + \text{gap} \\
 & && B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3 \geq a_2 + \text{gap} \\
 & && B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3 \geq a_1 + \text{gap} \\
 & && -2 \leq a \leq 2 \\
 & && -2 \leq B \leq 2
 \end{aligned} \tag{992}$$

having solution:

$$a^* = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}, \quad B^* = \begin{bmatrix} 0 & 2 & 2 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{bmatrix}, \quad \text{gap}^* = 1 \tag{993}$$

easily found by cvx [191] under MATLAB. For higher-dimensional  $q$  vectors (by induction),

$$a^* = -\mathbf{1} \in \mathbb{R}^n, \quad B^* = \begin{bmatrix} 0 & & \mathbf{2} \\ & \ddots & \\ \mathbf{0} & & 0 \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad \text{gap}^* = 1 \tag{994}$$

### 4.8.2 quantum Eternity II

Any equality of the form  $Ax = b$ , having binary solution  $x$ , may be expressed as a QUBO

$$\underset{x \in \{0, 1\}^n}{\text{minimize}} \quad x^T A^T A x - 2x^T A^T b \tag{995}$$

(§E.0.1.0.1) where  $B \triangleq A^T A - \delta^2(A^T A)$  and  $a \triangleq \delta(A^T A) - 2A^T b$  from (990). An adiabatic quantum annealer (like D:Wave's) is theoretically capable of solving Eternity II because it may be expressed  $\tilde{E}q = \tilde{r}$  (970) assuming that any feasible binary solution is minimal cardinality (p.317). This formulation (995) decreases sparsity, from that of  $A$ , which increases required qubit coupling.<sup>4.96</sup>

<sup>4.96</sup>For sparsity as defined on page 315, for nonsymmetric  $B$  matrix, and for:

- matrix  $E$  corresponding to (964), sparsity decreases from 0.0000052771 to 0.002683
- matrix  $\tilde{E}$  corresponding to (971), sparsity decreases from 0.00051786 to 0.027965
- matrix  $\tilde{E}$  corresponding to (972), sparsity decreases from 0.00056985 to 0.0047694
- matrix  $\tilde{E}$  corresponding to (977), sparsity decreases from 0.00070522 to 0.0042453.

**4.8.2.0.1 Example.** (E. D. Dahl) Nonincreasing discrete step.

		quantum step: $\dots$
$q_1$	$q_2$	$B_{12}q_1q_2 + a_1q_1 + a_2q_2$
desirable		
0	0	0
1	0	$a_1$
1	1	$B_{12} + a_1 + a_2$
undesirable		
0	1	$a_2$

$$\begin{aligned}
 & \underset{B, a, \text{gap}}{\text{maximize}} && \text{gap} \\
 & \text{subject to} && a_2 \geq 0 \quad + \text{gap} \\
 & && a_2 \geq a_1 \quad + \text{gap} \\
 & && a_2 \geq B_{12} + a_1 + a_2 \quad + \text{gap} \\
 & && -\mathbf{1} \leq a \leq \mathbf{1} \\
 & && -\mathbf{1} \leq B \leq \mathbf{1}
 \end{aligned} \tag{996}$$

Upper and lower bounds are 1, on each entrywise inequality, because gap is sufficient;

$$a^* = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \quad B^* = \begin{bmatrix} 0 & -1 \\ 0 & 0 \end{bmatrix}, \quad \text{gap}^* = 1 \tag{997}$$

Extensible to higher dimension; e.g.  $\{000, 100, 110, 111\}^T$  are desirable  $q \in \mathbb{R}^3$ .  $\square$

**4.8.2.0.2 Example.** (E. D. Dahl) Boolean qubit AND function.

We consider the case where second argument to AND is complemented:

			quantum AND function: $q_3 = q_1 \cdot \setminus q_2$
$q_1$	$q_2$	$q_3$	$B_{12}q_1q_2 + B_{13}q_1q_3 + B_{23}q_2q_3 + a_1q_1 + a_2q_2 + a_3q_3$
desirable			
0	0	0	0
0	1	0	$a_2$
1	0	1	$B_{13} + a_1 + a_3$
1	1	0	$B_{12} + a_1 + a_2$
undesirable			
0	0	1	$a_3$
0	1	1	$B_{23} + a_2 + a_3$
1	0	0	$a_1$
1	1	1	$B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3$

$$\begin{aligned}
& \underset{B, a, \text{gap}}{\text{maximize}} && \text{gap} \\
& \text{subject to} && a_3 \geq 0 \quad + \text{gap} \\
& && a_3 \geq a_2 \quad + \text{gap} \\
& && a_3 \geq B_{13} + a_1 + a_3 \quad + \text{gap} \\
& && a_3 \geq B_{12} + a_1 + a_2 \quad + \text{gap} \\
& && B_{23} + a_2 + a_3 \geq 0 \quad + \text{gap} \\
& && B_{23} + a_2 + a_3 \geq a_2 \quad + \text{gap} \\
& && B_{23} + a_2 + a_3 \geq B_{13} + a_1 + a_3 \quad + \text{gap} \\
& && B_{23} + a_2 + a_3 \geq B_{12} + a_1 + a_2 \quad + \text{gap} \\
& && a_1 \geq 0 \quad + \text{gap} \\
& && a_1 \geq a_2 \quad + \text{gap} \\
& && a_1 \geq B_{13} + a_1 + a_3 \quad + \text{gap} \\
& && a_1 \geq B_{12} + a_1 + a_2 \quad + \text{gap} \\
& && B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3 \geq 0 \quad + \text{gap} \\
& && B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3 \geq a_2 \quad + \text{gap} \\
& && B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3 \geq B_{13} + a_1 + a_3 \quad + \text{gap} \\
& && B_{12} + B_{13} + B_{23} + a_1 + a_2 + a_3 \geq B_{12} + a_1 + a_2 \quad + \text{gap} \\
& && -2 \leq a \leq 2 \\
& && -2 \leq B \leq 2
\end{aligned} \tag{998}$$

Optimal coefficients are not unique, but optimal objective gap is:

$$a^* = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}, \quad B^* = \begin{bmatrix} 0 & 1 & -2 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{bmatrix}, \quad \text{gap}^* = 1 \tag{999}$$

This optimal  $B$  matrix represents required coupling for AND but cannot be implemented in Chimera directly because there is no completely coupled three-qubit circuit.  $\square$

## 4.9 Constraining rank of indefinite matrices

Example 4.9.0.0.1, which follows, demonstrates that convex iteration is more generally applicable to indefinite or nonsquare matrices  $X \in \mathbb{R}^{m \times n}$ ; not only to positive semidefinite matrices. Indeed,

$$\begin{aligned}
& \underset{X \in \mathbb{R}^{m \times n}}{\text{find}} && X \\
& \text{subject to} && X \in \mathcal{C} \\
& && \text{rank } X \leq k
\end{aligned} \equiv \begin{aligned}
& \underset{X, Y, Z}{\text{find}} && X \\
& \text{subject to} && X \in \mathcal{C} \\
& && G = \begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \\
& && \text{rank } G \leq k
\end{aligned} \tag{1000}$$

**Proof.**  $\text{rank } G \leq k \Rightarrow \text{rank } X \leq k$  because  $X$  is the projection of composite matrix  $G$  on subspace  $\mathbb{R}^{m \times n}$ . For symmetric  $Y$  and  $Z$ , any rank- $k$  positive semidefinite composite matrix  $G$  can be factored into rank- $k$  terms  $R$ :  $G = R^T R$  where  $R \triangleq [B \ C]$  and  $\text{rank } B, \text{rank } C \leq \text{rank } R$  and  $B \in \mathbb{R}^{k \times m}$  and  $C \in \mathbb{R}^{k \times n}$ . Because  $Y$  and  $Z$  and  $X = B^T C$  are variable, (1619)  $\text{rank } X \leq \text{rank } B, \text{rank } C \leq \text{rank } R = \text{rank } G$  is tight.  $\spadesuit$

So, there must exist an optimal direction vector  $W^*$  such that

$$\begin{array}{ll} \text{find}_{X,Y,Z} & X \\ \text{subject to} & X \in \mathcal{C} \\ & G = \begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \\ & \text{rank } G \leq k \end{array} \quad \equiv \quad \begin{array}{ll} \text{minimize}_{X,Y,Z} & \langle G, W^* \rangle \\ \text{subject to} & X \in \mathcal{C} \\ & G = \begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \succeq 0 \end{array} \quad (1001)$$

Were  $W^* = I$ , the optimal resulting trace objective would be equivalent to the minimization of nuclear norm of  $X$  over  $\mathcal{C}$  by (1857). This means:

- (confer p.177) The argument of any nuclear norm minimization problem may be replaced with a composite semidefinite variable of the same optimal rank but doubly dimensioned.

Then Figure 96 becomes an accurate geometrical description of a consequent composite semidefinite problem objective. But there are better direction vectors than Identity  $I$  which occurs only under special conditions:

#### 4.9.0.0.1 Example. Compressed sensing, compressive sampling. [336]

*As our modern technology-driven civilization acquires and exploits ever-increasing amounts of data, everyone now knows that most of the data we acquire can be thrown away with almost no perceptual loss - witness the broad success of lossy compression formats for sounds, images, and specialized technical data. The phenomenon of ubiquitous compressibility raises very natural questions: Why go to so much effort to acquire all the data when most of what we get will be thrown away? Can't we just directly measure the part that won't end up being thrown away?*

–David Donoho [136]

Lossy data compression techniques like JPEG are popular, but it is also well known that compression artifacts become quite perceptible with signal postprocessing that goes beyond mere playback of a compressed signal. [247] [274] Spatial or audio frequencies presumed masked by a simultaneity are not encoded, for example, so rendered imperceptible even with significant postfiltering (of the compressed signal) that is meant to reveal them; *id est*, desirable artifacts are forever lost, so highly compressed data is not amenable to analysis and postprocessing: *e.g.*, sound effects [108] [109] [111] or image enhancement (Adobe Photoshop). <sup>4.97</sup> Further, there can be no universally acceptable unique metric of perception for gauging exactly how much data can be tossed. For these reasons, there will always be need for raw (noncompressed) data.

In this example, only so much information is thrown out as to leave perfect reconstruction within reach. Specifically, the MIT logo in Figure 135 is perfectly reconstructed from 700 time-sequential samples  $\{y_i\}$  acquired by the one-pixel camera illustrated in Figure 136. The MIT-logo image in this example impinges a  $46 \times 81$  array micromirror device. This mirror array is modulated by a pseudonoise source that independently positions all the individual mirrors. A single photodiode (one pixel) integrates incident light from all mirrors. After stabilizing the mirrors to a fixed but pseudorandom pattern, light so collected is then digitized into one sample  $y_i$  by analog-to-digital (A/D) conversion. This sampling process is repeated with the micromirror array modulated to a new pseudorandom pattern.

The most important questions are: How many samples are needed for perfect reconstruction? Does that number of samples represent compression of the original data?

---

<sup>4.97</sup> As simple a process as upward scaling of signal amplitude or image size will always introduce noise; even to a noncompressed signal. But scaling-noise is particularly noticeable in a JPEG-compressed image; *e.g.*, text or any sharp edge.



Figure 135: Massachusetts Institute of Technology (MIT) logo, including its white boundary, may be interpreted as a rank-5 matrix. This constitutes *Scene Y* observed by the one-pixel camera in Figure 136 for Example 4.9.0.0.1.

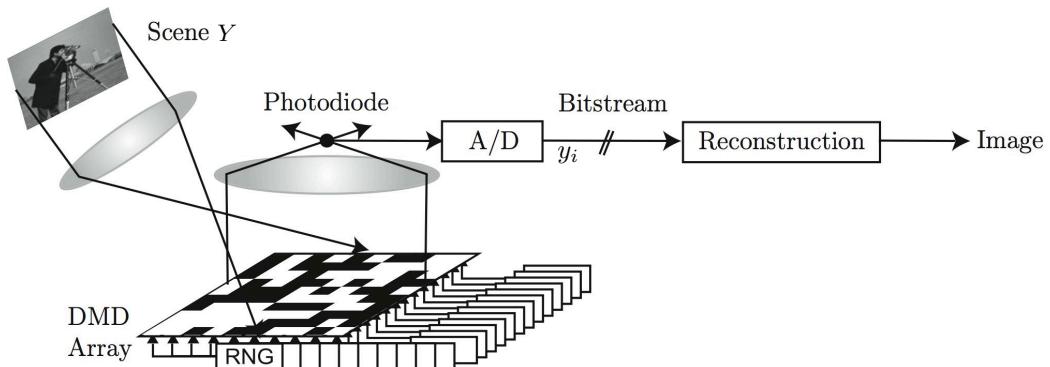


Figure 136: One-pixel camera. *Compressive imaging camera block diagram.* Incident lightfield (corresponding to the desired image  $Y$ ) is reflected off a digital micromirror device (DMD) array whose mirror orientations are modulated in the pseudorandom pattern supplied by the random number generators (RNG). Each different mirror pattern produces a voltage at the single photodiode that corresponds to one measurement  $y_i$ . –[383] [432]

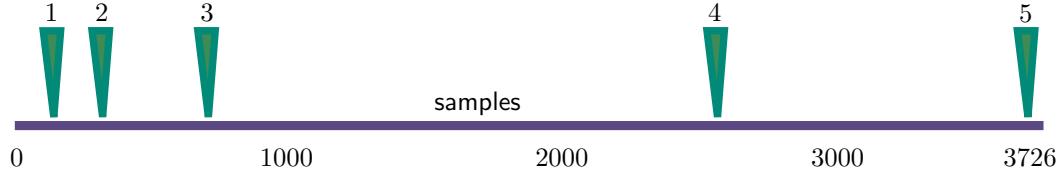


Figure 137: Estimates of compression for various encoding methods:

- 1) linear interpolation (140 samples),
- 2) minimal columnar basis (311 samples),
- 3) convex iteration (700 samples) can achieve lower bound predicted by compressed sensing (670 samples,  $n=46 \times 81$ ,  $k=140$ , Figure 113) whereas nuclear norm minimization alone does not [336, §6],
- 4) JPEG @100% quality (2588 samples),
- 5) no compression (3726 samples).

We claim that perfect reconstruction of the MIT logo can be achieved reliably with as few as 700 samples  $y = [y_i] \in \mathbb{R}^{700}$  from this one-pixel camera. That number represents only 19% of information obtainable from 3726 micromirrors.<sup>4.98</sup> (Figure 137)

Our approach to reconstruction is to look for low-rank solution to an *underdetermined* system:

$$\begin{array}{ll} \text{find} & X \\ X \in \mathbb{R}^{46 \times 81} & \\ \text{subject to} & A \text{ vec } X = y \\ & \text{rank } X \leq 5 \end{array} \quad (1002)$$

where  $\text{vec } X$  is the vectorized (37) unknown image matrix. Each row of wide matrix  $A$  is one realization of a pseudorandom pattern applied to the micromirrors. Since these patterns are deterministic (known), then the  $i^{\text{th}}$  sample  $y_i$  equals  $A(i, :) \text{ vec } Y$ ; *id est*,  $y = A \text{ vec } Y$ . *Perfect reconstruction* here means optimal solution  $X^*$  equals scene  $Y \in \mathbb{R}^{46 \times 81}$  to within machine precision.

Because variable matrix  $X$  is generally not square or positive semidefinite, we constrain its rank by rewriting the problem equivalently

$$\begin{array}{ll} \text{find} & X \\ W_1 \in \mathbb{R}^{46 \times 46}, W_2 \in \mathbb{R}^{81 \times 81}, X \in \mathbb{R}^{46 \times 81} & \\ \text{subject to} & A \text{ vec } X = y \\ & \text{rank} \begin{bmatrix} W_1 & X \\ X^T & W_2 \end{bmatrix} \leq 5 \end{array} \quad (1003)$$

This rank constraint on the composite (block) matrix insures  $\text{rank } X \leq 5$  for any choice of dimensionally compatible matrices  $W_1$  and  $W_2$ . But to solve this problem by convex iteration, we alternate solution of semidefinite program

$$\begin{array}{ll} \text{minimize}_{W_1 \in \mathbb{S}^{46}, W_2 \in \mathbb{S}^{81}, X \in \mathbb{R}^{46 \times 81}} & \text{tr} \left( \begin{bmatrix} W_1 & X \\ X^T & W_2 \end{bmatrix} Z \right) \\ \text{subject to} & A \text{ vec } X = y \\ & \begin{bmatrix} W_1 & X \\ X^T & W_2 \end{bmatrix} \succeq 0 \end{array} \quad (1004)$$

<sup>4.98</sup>That number (700 samples) is difficult to achieve, as reported in [336, §6]. If a minimal basis for the MIT logo were instead constructed, only five rows or columns worth of data (from a  $46 \times 81$  matrix) are linearly independent. This means a lower bound on achievable compression is about  $5 \times 46 = 230$  samples plus 81 samples column encoding; which corresponds to 8% of the original information. (Figure 137)

with semidefinite program

$$\begin{array}{ll} \text{minimize}_{Z \in \mathbb{S}^{46+81}} & \text{tr} \left( \begin{bmatrix} W_1 & X \\ X^T & W_2 \end{bmatrix}^* Z \right) \\ \text{subject to} & 0 \preceq Z \preceq I \\ & \text{tr } Z = 46 + 81 - 5 \end{array} \quad (1005)$$

(which has an optimal solution known in closed form, p.533) until a rank-5 composite matrix is found.

With 1000 samples  $\{y_i\}$ , convergence occurs in two iterations; 700 samples require more than ten iterations but reconstruction remains perfect. Iterating more admits taking of fewer samples. Reconstruction is independent of pseudorandom sequence parameters; *e.g.*, binary sequences succeed with the same efficiency as Gaussian or uniformly distributed sequences.  $\square$

#### 4.9.1 rank-constraint midsummary

We find that this *direction matrix* idea works well and quite independently of desired rank or affine dimension. This idea of direction matrix is good principally because of its simplicity: When confronted with a problem otherwise convex if not for a rank or cardinality constraint, then that constraint becomes a linear regularization term in the objective.

There exists a common thread through all these Examples; that being, convex iteration with a direction matrix as normal to a linear regularization (a generalization of the well-known trace heuristic). But each problem type (per Example) possesses its own idiosyncrasies that slightly modify how a rank-constrained optimal solution is actually obtained: The *ball packing* problem in Chapter 5.4.2.2.6, for example, requires a problem sequence in a progressively larger number of balls to find a good initial value for the direction matrix, whereas many of the examples in the present chapter require an initial value of  $\mathbf{0}$ . Finding a Boolean solution in Example 4.7.0.0.9 requires a procedure to detect stalls, while other problems have no such requirement. The combinatorial Procrustes problem in Example 4.7.0.0.3 allows use of a known closed-form solution for direction vector when solved via rank constraint, but not when solved via cardinality constraint. Some problems require a careful weighting of the regularization term, whereas other problems do not, and so on. It would be nice if there were a universally applicable method for constraining rank; one that is less susceptible to quirks of a particular problem type.

Poor initialization of the direction matrix from the regularization can lead to an erroneous result. We speculate one reason to be a simple dearth of optimal solutions of desired rank or cardinality,<sup>4.99</sup> an unfortunate choice of initial search direction leading astray. Ease of solution by convex iteration occurs when optimal solutions abound. With this speculation in mind, we now propose a further generalization of convex iteration for constraining rank that attempts to ameliorate quirks and unify problem types:

## 4.10 Convex Iteration rank-1

We now develop a general method for constraining rank that first decomposes a given problem via standard diagonalization of matrices (§A.5). This method is motivated by observation (§4.5.1.1) that an optimal direction matrix can be simultaneously diagonalizable with an optimal variable matrix. This suggests minimization of an

---

<sup>4.99</sup>In Convex Optimization, an optimal solution generally comes from a convex set of optimal solutions; (§3.1.1.1) that set can be large.

objective function directly in terms of eigenvalues. A second motivating observation is that variable orthogonal matrices seem easily found by convex iteration; *e.g.*, Procrustes Example 4.7.0.0.2.

### 4.10.1 rank-1 transformation

It turns out that this general method always requires solution to a rank-1 constrained problem regardless of desired rank  $\rho$  from the original problem. To demonstrate, we pose a semidefinite feasibility problem

$$\begin{aligned} \text{find } & X \in \mathbb{S}^n \\ \text{subject to } & A \text{ svec } X = b \\ & X \succeq 0 \\ & \text{rank } X \leq \rho \end{aligned} \tag{1006}$$

given an upper bound  $0 < \rho < n$  on rank, a vector  $b \in \mathbb{R}^m$ , and typically wide full-rank

$$A = \begin{bmatrix} \text{svec}(A_1)^T \\ \vdots \\ \text{svec}(A_m)^T \end{bmatrix} \in \mathbb{R}^{m \times n(n+1)/2} \tag{698}$$

where  $A_i \in \mathbb{S}^n$ ,  $i = 1 \dots m$ . So, for symmetric matrix vectorization svec as defined in (57),

$$A \text{ svec } X = \begin{bmatrix} \text{tr}(A_1 X) \\ \vdots \\ \text{tr}(A_m X) \end{bmatrix} \tag{699}$$

This program (1006) is a statement of the classical problem of finding a matrix  $X$  of maximum rank  $\rho$  in the intersection of the positive semidefinite cone with a given number  $m$  of hyperplanes in the subspace of symmetric matrices  $\mathbb{S}^n$ . [27, §II.13] [25, §2.2] Such a matrix is presumed to exist.

To begin transformation of (1006), express the nonincreasingly ordered diagonalization (§A.5.1) of positive semidefinite variable matrix

$$X \triangleq Q \Lambda Q^T = \sum_{i=1}^n \lambda_i Q_{ii} \in \mathbb{S}^n \tag{1007}$$

which is a sum of rank-1 orthogonal-projection matrices  $Q_{ii}$  weighted by eigenvalues  $\lambda_i$  where  $Q_{ij} \triangleq q_i q_j^T \in \mathbb{R}^{n \times n}$ ,  $Q = [q_1 \dots q_n] \in \mathbb{R}^{n \times n}$ ,  $Q^T = Q^{-1}$ ,  $\Lambda_{ii} = \lambda_i \in \mathbb{R}$ , and

$$\Lambda = \begin{bmatrix} \lambda_1 & & \mathbf{0} \\ & \lambda_2 & \\ & & \ddots \\ \mathbf{0}^T & & \lambda_n \end{bmatrix} \in \mathbb{S}^n \tag{1008}$$

where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ . Recall the fact:

$$\Lambda \succeq 0 \Leftrightarrow X \succeq 0 \tag{1604}$$

From orthogonal matrix  $Q$  in ordered diagonalization (1007) of variable  $X$ , take a matrix

$$U \triangleq [u_1 \dots u_\rho] \triangleq Q(:, 1:\rho) \sqrt{\Lambda(1:\rho, 1:\rho)} = [\sqrt{\lambda_1} q_1 \dots \sqrt{\lambda_\rho} q_\rho] \in \mathbb{R}^{n \times \rho} \tag{1009}$$

Then  $U$  has orthogonal but unnormalized columns;

$$X = UU^T = \sum_{i=1}^{\rho} u_i u_i^T \triangleq \sum_{i=1}^{\rho} U_{ii} = \sum_{i=1}^{\rho} \lambda_i q_i q_i^T \in \mathbb{S}^n \quad (1010)$$

Make an assignment

$$\begin{aligned} Z &= \begin{bmatrix} u_1 \\ \vdots \\ u_\rho \end{bmatrix} \begin{bmatrix} u_1^T & \cdots & u_\rho^T \end{bmatrix} \in \mathbb{S}^{n\rho} \\ &= \begin{bmatrix} U_{11} & \cdots & U_{1\rho} \\ \vdots & \ddots & \vdots \\ U_{1\rho}^T & \cdots & U_{\rho\rho} \end{bmatrix} \triangleq \begin{bmatrix} u_1 u_1^T & \cdots & u_1 u_\rho^T \\ \vdots & \ddots & \vdots \\ u_\rho u_1^T & \cdots & u_\rho u_\rho^T \end{bmatrix} \end{aligned} \quad (1011)$$

Then transformation of (1006) to its rank-1 equivalent is:

$$\begin{aligned} \text{find}_{U_{ii} \in \mathbb{S}^n, U_{ij} \in \mathbb{R}^{n \times n}} \quad X &= \sum_{i=1}^{\rho} U_{ii} \\ \text{subject to} \quad Z &= \begin{bmatrix} U_{11} & \cdots & U_{1\rho} \\ \vdots & \ddots & \vdots \\ U_{1\rho}^T & \cdots & U_{\rho\rho} \end{bmatrix} (\succeq 0) \\ A \operatorname{svec} \sum_{i=1}^{\rho} U_{ii} &= b \\ \operatorname{tr} U_{ij} &= 0 \quad i < j = 2 \dots \rho \\ \operatorname{rank} Z &= 1 \end{aligned} \quad (1012)$$

Symmetry is necessary and sufficient for positive semidefiniteness of a rank-1 matrix. (§A.3.1.0.7) Matrix  $X$  is positive semidefinite whenever  $Z$  is. (§A.3.1.0.4, §A.3.1.0.2) This new problem always enforces a rank-1 constraint on matrix  $Z$ ; *id est*, regardless of upper bound on rank  $\rho$  of variable matrix  $X$ , this equivalent problem always poses a rank-1 constraint. Upper bound  $\rho$  on rank of positive semidefinite matrix  $X$  is assured by rank-1 optimal matrix  $Z$ .

We propose solving (1012) by iteration of convex problem

$$\begin{aligned} \text{minimize}_{U_{ii} \in \mathbb{S}^n, U_{ij} \in \mathbb{R}^{n \times n}} \quad &\operatorname{tr}(Z W) \\ \text{subject to} \quad Z &= \begin{bmatrix} U_{11} & \cdots & U_{1\rho} \\ \vdots & \ddots & \vdots \\ U_{1\rho}^T & \cdots & U_{\rho\rho} \end{bmatrix} \succeq 0 \\ A \operatorname{svec} \sum_{i=1}^{\rho} U_{ii} &= b \\ \operatorname{tr} U_{ij} &= 0 \quad i < j = 2 \dots \rho \end{aligned} \quad (1013)$$

with convex problem

$$\begin{aligned} \text{minimize}_{W \in \mathbb{S}^{n\rho}} \quad &\operatorname{tr}(Z^* W) \\ \text{subject to} \quad 0 &\preceq W \preceq I \\ \operatorname{tr} W &= n\rho - 1 \end{aligned} \quad (1014)$$

the latter providing direction of search  $W$  for a rank-1 matrix  $Z$  in (1013). These convex problems (1013) (1014) are iterated until a rank-1  $Z$  matrix is found (until the objective

of (1013) vanishes). Initial value of direction matrix  $W$  is the Identity. For subsequent iterations, an optimal solution to (1014) has closed form (p.533).

Because of the nonconvex nature of a rank-constrained problem, there can be no proof of convergence of this convex iteration to a feasible point of (1012). But the iteration always converges to a local minimum because the sequence of objective values is monotonic and nonincreasing; any monotonically nonincreasing real sequence converges. [289, §1.2] [43, §1.1] A rank  $\rho$  matrix  $X$  solving the original problem (1006) is found when the objective in (1013) converges to 0: a certificate of global optimality for the convex iteration. In practice, incidence of success is quite high (99.99% [414]); failures being mostly attributable to numerical accuracy.

#### 4.10.2 singular value decomposition by convex iteration

This diagonal decomposition technique (transformation to a rank-1 problem) is extensible to other problem types; *e.g.*, [252, §III]. Rank-1 transformation makes singular value decomposition (SVD, §A.6) possible by convex iteration because orthogonality constraints may then be introduced. We learn that any uniqueness properties, the SVD of rank- $\rho$  matrix

$$X \triangleq USV^T \in \mathbb{R}^{m \times n} \quad (1015)$$

might enjoy, stem from demand for singular vector orthonormality.<sup>4.100</sup>

Assignment  $Z \in \mathbb{S}_+^{2m\rho+n\rho+\rho+1}$  is key to finding the SVD of  $X$  by convex optimization:

$$\begin{aligned} & \underset{H, J}{\text{find}} \quad U, \delta(S), V \\ & \text{subject to} \quad Z = \begin{bmatrix} 1 & \text{vec}(H)^T & \text{vec}(U)^T & \delta(S)^T & \text{vec}(V)^T \\ \text{vec } H & & & & \\ \text{vec } U & & J & & \\ \delta(S) & & & & \\ \text{vec } V & & & & \end{bmatrix} \succeq 0 \\ & \delta(S) \succeq 0 \\ & H = US \subset J \\ & X = HV^T \in J \\ & HU^T \text{ symmetry} \\ & U^T H \text{ perpendicularity} \\ & \text{tr}(H(:, i)H(:, i)^T) = S(i, i)^2 \quad i=1 \dots \rho \\ & \text{tr}(H(:, i)U(:, i)^T) = S(i, i) \quad i=1 \dots \rho \\ & H \text{ orthogonality} \\ & U \text{ orthonormality} \\ & V \text{ orthonormality} \\ & \text{rank } Z = 1 \end{aligned} \quad (1016)$$

where variable matrix  $J \in \mathbb{S}_+^{2m\rho+n\rho+\rho}$  is a large partition of  $Z$ , where given rank- $\rho$  matrix  $X \in \mathbb{R}^{m \times n}$  is subject to SVD in unknown orthonormal matrices  $U \in \mathbb{R}^{m \times \rho}$  and  $V \in \mathbb{R}^{n \times \rho}$  and unknown diagonal matrix of singular values  $S \in \mathbb{R}^{\rho \times \rho}$ , and where introduction of variable

$$H \triangleq US \in \mathbb{R}^{m \times \rho} \quad (1017)$$

makes identification of input  $X = HV^T$  possible within partition  $J$ . Orthogonality constraints on columns of  $H$ , within  $J$ , and orthonormality constraints on columns of  $U$  and  $V$  are critical; *videlicet*,  $h \perp v \Leftrightarrow \text{tr}(hv^T)=0$ ;  $v^T v = 1 \Leftrightarrow \text{tr}(vv^T)=1$ .

---

<sup>4.100</sup>Otherwise, there exist many similarly structured tripartite nonorthogonal matrix decompositions; in place of  $\rho$  nonzero singular values, diagonal matrix  $S$  would instead hold exactly  $\rho$  coordinates; orthonormal columns in  $U$  and  $V$  would become merely linearly independent.

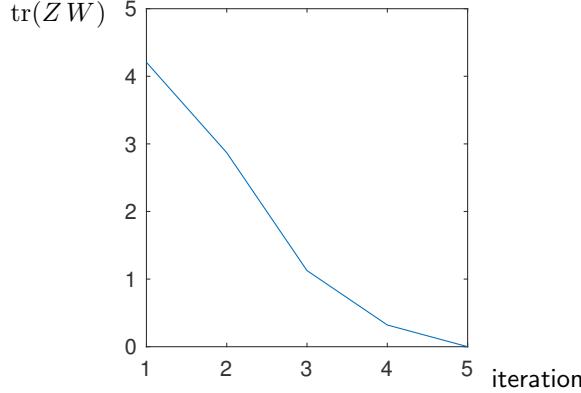


Figure 138: Typical convergence of SVD by convex iteration for a  $2 \times 2$  random  $X$  matrix. Matrix  $W$  represents a direction vector of convex iteration rank-1.

Symmetric matrix  $Z$  is positive semidefinite rank-1 at optimality, regardless of rank  $\rho$ . That rank constraint is the only nonconvex constraint in (1016); the only constraint that cannot be directly implemented in a convex manner per partition  $J$ . But the rank constraint is handled well by convex iteration. MATLAB implementation of SVD by convex iteration is intricate although incidence of success is 99.99% [429], barring numerical error.

#### 4.10.2.0.1 Example. SVD of $X$ by convex iteration.

Given rank-2 matrix  $X = USV^T \in \mathbb{R}^{2 \times 2}$ , we now make explicit every constraint in (1016):

$$\begin{aligned} & \text{find}_{H \in \mathbb{R}^{2 \times 2}, J \in \mathbb{S}^{14}} \quad U \in \mathbb{R}^{2 \times 2}, S = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}, V \in \mathbb{R}^{2 \times 2} \\ & \text{subject to} \quad Z = \begin{bmatrix} 1 & h_1^T & h_2^T & u_1^T & u_2^T & [\sigma_1 \ \sigma_2] & v_1^T & v_2^T \\ h_1 & J_{11} & J_{12} & J_{13} & J_{14} & J_{15} & J_{16} & J_{17} \\ h_2 & J_{12}^T & J_{22} & J_{23} & J_{24} & J_{25} & J_{26} & J_{27} \\ u_1 & J_{13}^T & J_{23}^T & J_{33} & J_{34} & J_{35} & J_{36} & J_{37} \\ u_2 & J_{14}^T & J_{24}^T & J_{34} & J_{44} & J_{45} & J_{46} & J_{47} \\ \begin{bmatrix} \sigma_1 \\ \sigma_2 \end{bmatrix} & J_{15}^T & J_{25}^T & J_{35}^T & J_{45}^T & J_{55} & J_{56} & J_{57} \\ v_1 & J_{16}^T & J_{26}^T & J_{36}^T & J_{46}^T & J_{56}^T & J_{66} & J_{67} \\ v_2 & J_{17}^T & J_{27}^T & J_{37}^T & J_{47}^T & J_{57}^T & J_{67}^T & J_{77} \end{bmatrix} \succeq 0 \end{aligned} \tag{1018}$$

$$\sigma_1, \sigma_2 \geq 0$$

$$H = [J_{35}(:, 1) \ J_{45}(:, 2)]$$

$$X = J_{16} + J_{27}$$

$$J_{13} = J_{13}^T, \ J_{24} = J_{24}^T$$

$$\text{tr } J_{14} = 0, \ \text{tr } J_{23} = 0$$

$$\text{tr } J_{11} = J_{55}(1, 1), \ \text{tr } J_{22} = J_{55}(2, 2)$$

$$\text{tr } J_{13} = \sigma_1, \ \text{tr } J_{24} = \sigma_2$$

$$\text{tr } J_{12} = 0$$

$$\text{tr } J_{33} = 1, \ \text{tr } J_{44} = 1, \ \text{tr } J_{34} = 0$$

$$\text{tr } J_{66} = 1, \ \text{tr } J_{77} = 1, \ \text{tr } J_{67} = 0$$

$$\text{rank } Z = 1$$

where  $H \triangleq [h_1 \ h_2]$ ,  $U \triangleq [u_1 \ u_2]$ ,  $S \triangleq \delta([\sigma_1 \ \sigma_2]^T)$ , and  $V \triangleq [v_1 \ v_2]$ . Observe how, excepting the rank constraint, constraints are written as affine expressions of variable

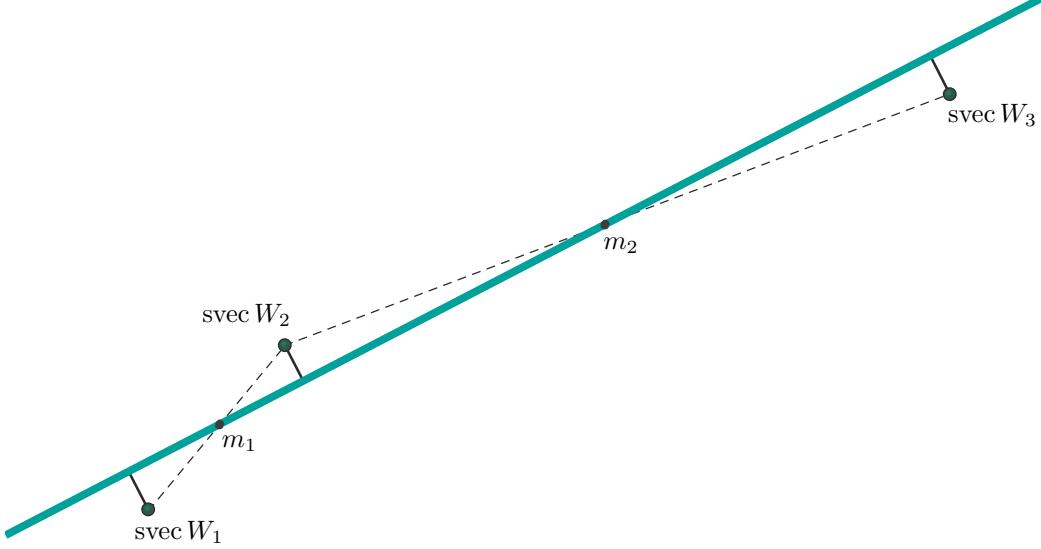


Figure 139:  $W_1$ ,  $W_2$ , and  $W_3$  represent the last three direction vectors in a sequence.  $m_1$  represents the midpoint between direction vectors  $W_1$  and  $W_2$ ;  $m_2$  is the midpoint of  $W_2$  and  $W_3$ . Straight line passes through midpoints.

matrix  $J$ . [429] Convergence is illustrated in Figure 138. Incidentally, the singular value decomposition of  $X = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$  has closed form (*Mathematica*)

$$USV^T =$$

$$\left[ \begin{array}{cc} \frac{a^3+2bcd+a(b^2+c^2-d^2-\gamma)}{\sqrt{(a^3+2bcd+a(b^2+c^2-d^2-\gamma))^2+(a^2c+2abd+c(-b^2+c^2+d^2-\gamma))^2}\text{sgn}(ab+cd)} & \frac{a^3+2bcd+a(b^2+c^2-d^2+\gamma)}{\sqrt{(a^3+2bcd+a(b^2+c^2-d^2+\gamma))^2+(a^2c+2abd+c(-b^2+c^2+d^2+\gamma))^2}\text{sgn}(ab+cd)} \\ \frac{a^2c+2abd+c(-b^2+c^2+d^2-\gamma)}{\sqrt{(a^3+2bcd+a(b^2+c^2-d^2-\gamma))^2+(a^2c+2abd+c(-b^2+c^2+d^2-\gamma))^2}\text{sgn}(ab+cd)} & \frac{a^2c+2abd+c(-b^2+c^2+d^2+\gamma)}{\sqrt{(a^3+2bcd+a(b^2+c^2-d^2+\gamma))^2+(a^2c+2abd+c(-b^2+c^2+d^2+\gamma))^2}\text{sgn}(ab+cd)} \\ \end{array} \right]$$

$$\left[ \begin{array}{cc} \frac{\sqrt{a^2+b^2+c^2+d^2-\gamma}}{\sqrt{2}} & 0 \\ 0 & \frac{\sqrt{a^2+b^2+c^2+d^2+\gamma}}{\sqrt{2}} \end{array} \right] \quad (1019)$$

$$\left[ \begin{array}{cc} \frac{a^2-b^2+c^2-d^2-\gamma}{(ab+cd)\sqrt{4+\frac{(-a^2+b^2-c^2+d^2+\gamma)^2}{(ab+cd)^2}}} & \frac{2}{\sqrt{4+\frac{(-a^2+b^2-c^2+d^2+\gamma)^2}{(ab+cd)^2}}} \\ \frac{a^2-b^2+c^2-d^2+\gamma}{(ab+cd)\sqrt{4+\frac{(a^2-b^2+c^2-d^2+\gamma)^2}{(ab+cd)^2}}} & \frac{2}{\sqrt{4+\frac{(a^2-b^2+c^2-d^2+\gamma)^2}{(ab+cd)^2}}} \end{array} \right]$$

where

$$\gamma \triangleq \sqrt{((b+c)^2 + (a-d)^2)((b-c)^2 + (a+d)^2)} \quad (1020)$$

□

## 4.11 Convex Iteration accelerant

Convex iteration can be made to converge faster; sometimes, by orders of magnitude. The idea here is to determine whether the last three direction vectors are close to their fit to a straight line. When three direction vectors are close to a straight line, then the last direction vector may be replaced with its extrapolation along that line.

To reduce computation time, a fitted line is not a best fit. Instead, the midpoint between each pair of iteration-adjacent direction vectors is calculated (Figure 139). A straight line is uniquely defined by two midpoints in any dimension. Distance of each direction vector to the line is calculated, then those three distances summed into a program variable called `straight`. When a sum is small, three direction vectors are deemed close to the line determined by them. What is meant by *close* and *small* depends on problem type and data. For the parameters and normalized random data chosen for two MATLAB realizations [414] [429] on *Wikimization* (corresponding to problems (1012) and (1016)), *small* is numerically defined to be 1 or less in the statement `if straight < 1` whose purpose is to determine straightness of the last three direction vectors of convex iteration. The smaller the value of sum `straight`, the closer the last three direction vectors are to a straight line. Variable `straight` is inherently bounded below by 0 which indicates three direction vectors precisely on the line going through them.

If linear extrapolation goes too far, then the objective of convex iteration will increase or a solver may fail numerically. In either case, one must forget the last iteration and back up the linear extrapolation until the objective decreases. These techniques are illustrated by the MATLAB programs; [414] Figure 138 is one representative. [429]



# Chapter 5

## Euclidean Distance Matrix

*These results [(1052)] were obtained by Schoenberg (1935), a surprisingly late date for such a fundamental property of Euclidean geometry.*

— John Clifford Gower [186, §3]

By itself, distance information between many points in Euclidean space is lacking. We might want to know more; such as, relative or absolute position or dimension of some hull. A question naturally arising in some fields (*e.g.*, geodesy, economics, genetics, psychology, biochemistry, engineering) [116] asks what facts can be deduced given only distance information. What can we know about the underlying points that the distance information purports to describe? We also ask what it means when given distance information is incomplete; or suppose the distance information is not reliable, available, or specified only by certain tolerances (affine inequalities). These questions motivate a study of interpoint distance, well represented in any spatial dimension by a simple matrix from linear algebra.<sup>5.1</sup> In what follows, we will answer some of these questions via Euclidean distance matrices.

### 5.1 EDM

Euclidean space  $\mathbb{R}^n$  is a finite-dimensional real vector space having an inner product defined on it, inducing a *metric*. [254, §3.1] A Euclidean distance matrix, an EDM in  $\mathbb{R}_+^{N \times N}$ , is an exhaustive table of distance-square  $d_{ij}$  between points taken by pair from a list of  $N$  points  $\{x_\ell, \ell=1 \dots N\}$  in  $\mathbb{R}^n$ ; the squared metric, the measure of distance-square:

$$d_{ij} = \|x_i - x_j\|_2^2 \triangleq \langle x_i - x_j, x_i - x_j \rangle \quad (1021)$$

Each point is labelled ordinally, hence the row or column index of an EDM,  $i$  or  $j=1 \dots N$ , individually addresses all the points in the list.

Consider the following example of an EDM for the case  $N=3$ :

---

<sup>5.1</sup> *e.g.*,  $\sqrt[D]{D} \in \mathbb{R}^{N \times N}$ , a classical two-dimensional matrix representation of absolute interpoint distance because its entries (in ordered rows and columns) can be written neatly on a piece of paper. Matrix  $D$  will be reserved throughout to hold distance-square.

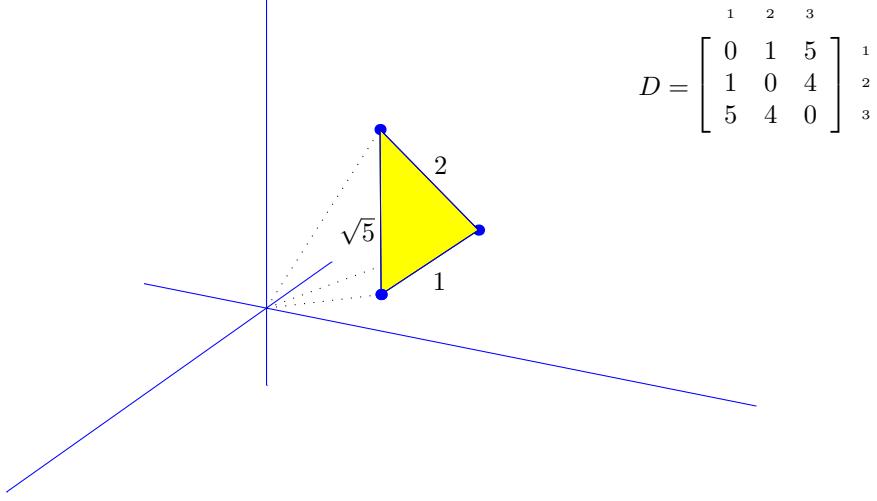


Figure 140: Convex hull of three points ( $N=3$ ) is shaded in  $\mathbb{R}^3$  ( $n=3$ ). Dotted lines are imagined vectors to points whose affine dimension is 2.

$$D = [d_{ij}] = \begin{bmatrix} d_{11} & d_{12} & d_{13} \\ d_{21} & d_{22} & d_{23} \\ d_{31} & d_{32} & d_{33} \end{bmatrix} = \begin{bmatrix} 0 & d_{12} & d_{13} \\ d_{12} & 0 & d_{23} \\ d_{13} & d_{23} & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 5 \\ 1 & 0 & 4 \\ 5 & 4 & 0 \end{bmatrix} \quad (1022)$$

Matrix  $D$  has  $N^2$  entries but only  $N(N-1)/2$  pieces of information. In Figure 140 are shown three points in  $\mathbb{R}^3$  that can be arranged in a list to correspond to  $D$  in (1022). But such a list is not unique because any rotation, reflection, or translation (§5.5) of those points would produce the same EDM  $D$ .

## 5.2 First metric properties

For  $i, j = 1 \dots N$ , absolute distance between points  $x_i$  and  $x_j$  must satisfy the defining requirements imposed upon any *metric space*: [254, §1.1] [289, §1.7] namely, for Euclidean metric  $\sqrt{d_{ij}}$  (§5.4) in  $\mathbb{R}^n$

- 1.  $\sqrt{d_{ij}} \geq 0$ ,  $i \neq j$  nonnegativity
- 2.  $\sqrt{d_{ii}} = 0 \Leftrightarrow x_i = x_j$  selfdistance
- 3.  $\sqrt{d_{ij}} = \sqrt{d_{ji}}$  symmetry
- 4.  $\sqrt{d_{ij}} \leq \sqrt{d_{ik}} + \sqrt{d_{kj}}$ ,  $i \neq j \neq k$  triangle inequality

Then all entries of an EDM must be in concord with these Euclidean metric properties: specifically, each entry must be nonnegative,<sup>5.2</sup> the main diagonal must be  $\mathbf{0}$ ,<sup>5.3</sup> and an EDM must be symmetric. The fourth property provides upper and lower bounds for each entry. Property 4 is true more generally when there are no restrictions on indices  $i, j, k$ , but furnishes no new information.

<sup>5.2</sup>Implicit from the terminology,  $\sqrt{d_{ij}} \geq 0 \Leftrightarrow d_{ij} \geq 0$  is always assumed.

<sup>5.3</sup>What we call selfdistance, Marsden calls *nondegeneracy*. [289, §1.6] Kreyszig calls these first metric properties *axioms of the metric*; [254, p.4] Blumenthal refers to them as *postulates*. [54, p.15]

### 5.3 $\exists$ fifth Euclidean metric property

The four properties of the Euclidean metric provide information insufficient to certify that a bounded convex polyhedron more complicated than a triangle has a Euclidean realization. [186, §2] Yet any list of points or the vertices of any bounded convex polyhedron must conform to the properties.

#### 5.3.0.0.1 Example. Triangle.

Consider the EDM in (1022), but missing one of its entries:

$$D = \begin{bmatrix} 0 & 1 & d_{13} \\ 1 & 0 & 4 \\ d_{31} & 4 & 0 \end{bmatrix} \quad (1023)$$

Can we determine unknown entries of  $D$  by applying the metric properties? Property 1 demands  $\sqrt{d_{13}}, \sqrt{d_{31}} \geq 0$ , property 2 requires the main diagonal be  $\mathbf{0}$ , while property 3 makes  $\sqrt{d_{31}} = \sqrt{d_{13}}$ . The fourth property tells us

$$1 \leq \sqrt{d_{13}} \leq 3 \quad (1024)$$

Indeed, described over that closed interval  $[1, 3]$  is a family of triangular polyhedra whose angle at vertex  $x_2$  varies from 0 to  $\pi$  radians. So, yes we can determine the unknown entries of  $D$ , but they are not unique; nor should they be from the information given for this example.  $\square$

#### 5.3.0.0.2 Example. Small completion problem, I.

Now consider the polyhedron in Figure 141b formed from an unknown list  $\{x_1, x_2, x_3, x_4\}$ . The corresponding EDM less one critical piece of information,  $d_{14}$ , is given by

$$D = \begin{bmatrix} 0 & 1 & 5 & d_{14} \\ 1 & 0 & 4 & 1 \\ 5 & 4 & 0 & 1 \\ d_{14} & 1 & 1 & 0 \end{bmatrix} \quad (1025)$$

From metric property 4 we may write a few inequalities for the two triangles common to  $d_{14}$ ; we find

$$\sqrt{5}-1 \leq \sqrt{d_{14}} \leq 2 \quad (1026)$$

We cannot further narrow those bounds on  $\sqrt{d_{14}}$  using only the four metric properties (§5.8.3.1.1). Yet there is only one possible choice for  $\sqrt{d_{14}}$  because points  $x_2, x_3, x_4$  must be collinear. All other values of  $\sqrt{d_{14}}$  in the interval  $[\sqrt{5}-1, 2]$  specify impossible distances in any dimension; *id est*, in this particular example the triangle inequality does not yield an interval for  $\sqrt{d_{14}}$  over which a family of convex polyhedra can be reconstructed.  $\square$

We will return to this simple Example 5.3.0.0.2 to illustrate more elegant methods of solution in §5.8.3.1.1, §5.9.3.0.1, and §5.14.4.1.1. Until then, we can deduce some general principles from the foregoing examples:

- Unknown  $d_{ij}$  of an EDM are not necessarily uniquely determinable.
- The triangle inequality does not produce necessarily tight bounds.<sup>5.4</sup>
- Four Euclidean metric properties are insufficient for reconstruction.

---

<sup>5.4</sup>The term *tight* with reference to an inequality means equality is achievable.

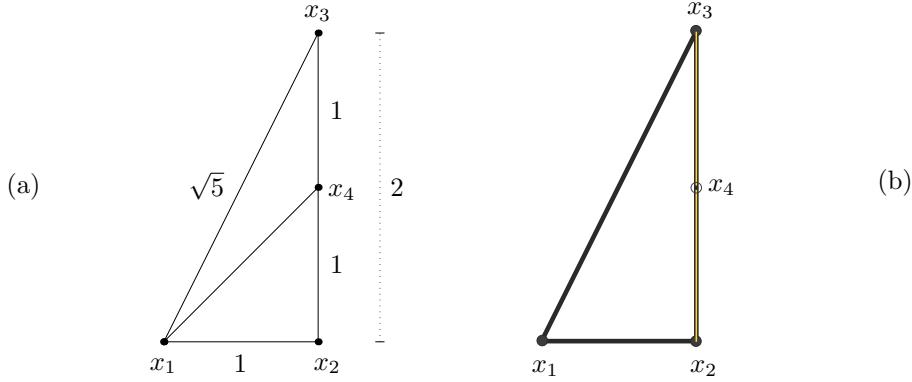


Figure 141: (a) Complete dimensionless *EDM graph*. (b) Emphasizing obscured segments  $\overline{x_2x_4}$ ,  $\overline{x_4x_3}$ , and  $\overline{x_2x_3}$ , now only five ( $2N-3$ ) absolute distances are specified. EDM so represented is incomplete, missing  $d_{14}$  as in (1025), yet the isometric reconstruction (§5.4.2.2.10) is unique as proved in §5.9.3.0.1 and §5.14.4.1.1. First four properties of Euclidean metric are not a recipe for reconstruction of this polyhedron.

### 5.3.1 lookahead

There must exist at least one requirement more than the four properties of the Euclidean metric that makes them altogether necessary and sufficient to certify realizability of bounded convex polyhedra. Indeed, there are infinitely many more; there are precisely  $N+1$  necessary and sufficient Euclidean metric requirements for  $N$  points constituting a generating list (§2.3.2). Here is the fifth requirement:

#### 5.3.1.0.1 Fifth Euclidean metric property. *Relative-angle inequality.*

(confer §5.14.2.1.1) Augmenting the four fundamental properties of the Euclidean metric in  $\mathbb{R}^n$ , for all  $i, j, \ell \neq k \in \{1 \dots N\}$ ,  $i < j < \ell$ , and for  $N \geq 4$  distinct points  $\{x_k\}$ , the inequalities

$$\begin{aligned} \cos(\theta_{ik\ell} + \theta_{\ellkj}) &\leq \cos \theta_{ikj} \leq \cos(\theta_{ik\ell} - \theta_{\ellkj}) \\ 0 \leq \theta_{ik\ell}, \theta_{\ellkj}, \theta_{ikj} &\leq \pi \end{aligned} \tag{1027}$$

where  $\theta_{ikj} = \theta_{jki}$  represents angle between vectors at vertex  $x_k$  (1099) (Figure 142), must be satisfied at each point  $x_k$  regardless of affine dimension.  $\diamond$

We will explore this in §5.14. One of our early goals is to determine matrix criteria that subsume all the Euclidean metric properties and any further requirements. Looking ahead, we will find (1378) (1052) (1056)

$$\left. \begin{array}{l} -z^T D z \geq 0 \\ \mathbf{1}^T z = 0 \\ (\forall \|z\| = 1) \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow D \in \mathbb{EDM}^N \tag{1028}$$

where the convex cone of Euclidean distance matrices  $\mathbb{EDM}^N \subseteq \mathbb{S}_h^N$  belongs to the subspace of symmetric hollow<sup>5.5</sup> matrices (§2.2.3.0.1). (Numerical test `isedm()` is provided on

<sup>5.5</sup> 0 main diagonal.

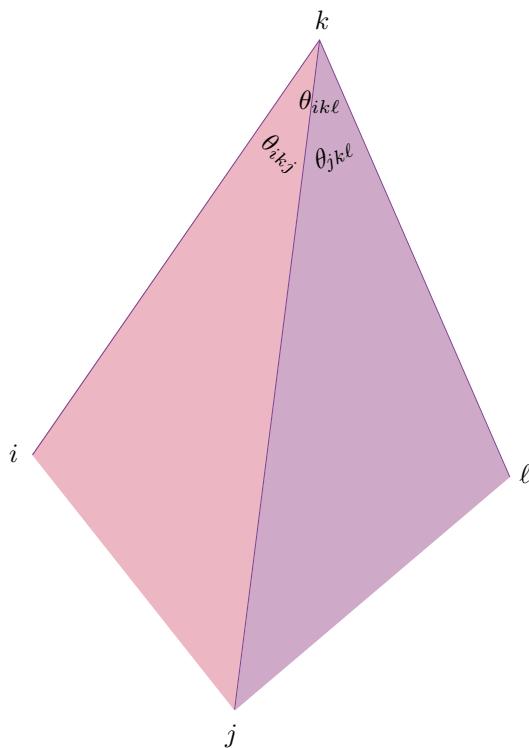


Figure 142: Fifth Euclidean metric property nomenclature. Each angle  $\theta$  is made by a vector pair at vertex  $k$  while  $i, j, k, \ell$  index four points at the vertices of a generally irregular tetrahedron. The fifth property is necessary for realization of four or more points; a reckoning by three angles in any dimension. Together with the first four Euclidean metric properties, this fifth property is necessary and sufficient for realization of four points.

*Wikimization* [431].) Having found equivalent matrix criteria, we will see there is a bridge from bounded convex polyhedra to EDMs in §5.9.<sup>5,6</sup>

Now we develop some invaluable concepts, moving toward a link of the Euclidean metric properties to matrix criteria.

## 5.4 EDM definition

Ascribe points in a list  $\{x_\ell \in \mathbb{R}^n, \ell = 1 \dots N\}$  to the columns of a matrix

$$X = [x_1 \ \cdots \ x_N] \in \mathbb{R}^{n \times N} \quad (77)$$

where  $N$  is regarded as *cardinality* of list  $X$ . When matrix  $D = [d_{ij}]$  is an EDM, its entries must be related to those points constituting the list by the Euclidean distance-square: for  $i, j = 1 \dots N$  (§A.1.1 no.36)

$$\begin{aligned} d_{ij} &= \|x_i - x_j\|^2 = (x_i - x_j)^T(x_i - x_j) = \|x_i\|^2 + \|x_j\|^2 - 2x_i^T x_j \\ &= [x_i^T \ x_j^T] \begin{bmatrix} I & -I \\ -I & I \end{bmatrix} \begin{bmatrix} x_i \\ x_j \end{bmatrix} \\ &= \text{vec}(X)^T(\Phi_{ij} \otimes I) \text{vec } X = \langle \Phi_{ij}, X^T X \rangle \end{aligned} \quad (1029)$$

where

$$\text{vec } X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^{nN} \quad (1030)$$

and where  $\otimes$  signifies Kronecker product (§D.1.2.1).  $\Phi_{ij} \otimes I$  is positive semidefinite (1636) having  $I \in \mathbb{S}^n$  in its  $ii^{\text{th}}$  and  $jj^{\text{th}}$  block of entries while  $-I \in \mathbb{S}^n$  fills its  $ij^{\text{th}}$  and  $ji^{\text{th}}$  block; *id est*,

$$\begin{aligned} \Phi_{ij} &\triangleq \delta((e_i e_j^T + e_j e_i^T) \mathbf{1}) - (e_i e_j^T + e_j e_i^T) \in \mathbb{S}_+^N \\ &= e_i e_i^T + e_j e_j^T - e_i e_j^T - e_j e_i^T \\ &= (e_i - e_j)(e_i - e_j)^T \end{aligned} \quad (1031)$$

where  $\{e_i \in \mathbb{R}^N, i = 1 \dots N\}$  is the set of standard basis vectors. Thus each entry  $d_{ij}$  is a convex quadratic function (§A.4.0.0.2) of  $\text{vec } X$  (37). [343, §6]

The collection of all Euclidean distance matrices  $\mathbb{EDM}^N$  is a convex subset of  $\mathbb{R}_+^{N \times N}$  called the *EDM cone* (§6, Figure 177 p.457);

$$\mathbf{0} \in \mathbb{EDM}^N \subseteq \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N} \subset \mathbb{S}^N \quad (1032)$$

An EDM  $D$  must be expressible as a function of some list  $X$ ; *id est*, it must have the form

$$\mathbf{D}(X) \triangleq \delta(X^T X) \mathbf{1}^T + \mathbf{1} \delta(X^T X)^T - 2X^T X \in \mathbb{EDM}^N \quad (1033)$$

$$= [\text{vec}(X)^T(\Phi_{ij} \otimes I) \text{vec } X, i, j = 1 \dots N] \quad (1034)$$

Function  $\mathbf{D}(X)$  will make an EDM given any  $X \in \mathbb{R}^{n \times N}$ , conversely, but  $\mathbf{D}(X)$  is not a convex function of  $X$  (§5.4.1). Now the EDM cone may be described:

$$\mathbb{EDM}^N = \left\{ \mathbf{D}(X) \mid X \in \mathbb{R}^{N-1 \times N} \right\} \quad (1035)$$

---

<sup>5,6</sup>From an EDM, a generating list (§2.3.2, §2.12.2) for a polyhedron can be found (§5.12) correct to within a rotation, reflection, and translation (§5.5).

Expression  $\mathbf{D}(X)$  is a matrix definition of EDM and so conforms to the Euclidean metric properties:

Nonnegativity of EDM entries (property 1, §5.2) is obvious from the distance-square definition (1029), so holds for any  $D$  expressible in the form  $\mathbf{D}(X)$  in (1033).

When we say  $D$  is an EDM, reading from (1033), it implicitly means the main diagonal must be  $\mathbf{0}$  (property 2, selfdistance) and  $D$  must be symmetric (property 3);  $\delta(D) = \mathbf{0}$  and  $D^T = D$  or, equivalently,  $D \in \mathbb{S}_h^N$  are necessary matrix criteria.

#### 5.4.0.1 homogeneity

Function  $\mathbf{D}(X)$  is homogeneous in the sense, for  $\zeta \in \mathbb{R}$

$$\sqrt[\nu]{\mathbf{D}(\zeta X)} = |\zeta| \sqrt[\nu]{\mathbf{D}(X)} \quad (1036)$$

where the positive square root is entrywise ( $\circ$ ).

Any nonnegatively scaled EDM remains an EDM; *id est*, the matrix class EDM is invariant to nonnegative scaling ( $\alpha \mathbf{D}(X)$  for  $\alpha \geq 0$ ) because all EDMs of dimension  $N$  constitute a convex cone  $\text{EDM}^N$  (§6, Figure 169).

#### 5.4.1 $-V_N^T \mathbf{D}(X) V_N$ convexity

We saw that EDM entries  $d_{ij} \left( \begin{bmatrix} x_i \\ x_j \end{bmatrix} \right)$  are convex quadratic functions. Yet  $-\mathbf{D}(X)$  (1033) is not a quasiconvex function of matrix  $X \in \mathbb{R}^{n \times N}$  because the second directional derivative (§3.14)

$$-\frac{d^2}{dt^2} \Big|_{t=0} \mathbf{D}(X + tY) = 2(-\delta(Y^T Y) \mathbf{1}^T - \mathbf{1} \delta(Y^T Y)^T + 2Y^T Y) \quad (1037)$$

is indefinite for any  $Y \in \mathbb{R}^{n \times N}$  since its main diagonal is  $\mathbf{0}$ . [181, §4.2.8] [228, §7.1 prob.2] Hence  $-\mathbf{D}(X)$  can neither be convex in  $X$ .

The outcome is different when instead we consider

$$-V_N^T \mathbf{D}(X) V_N = 2V_N^T X^T X V_N \quad (1038)$$

where we introduce the full-rank thin Schoenberg auxiliary matrix (§B.4.2)

$$V_N \triangleq \frac{1}{\sqrt{2}} \begin{bmatrix} -1 & -1 & \cdots & -1 \\ 1 & & & \mathbf{0} \\ & 1 & & \\ & & \ddots & \\ \mathbf{0} & & & 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} -\mathbf{1}^T \\ I \end{bmatrix} \in \mathbb{R}^{N \times N-1} \quad (1039)$$

$(\mathcal{N}(V_N) = \mathbf{0})$  having range

$$\mathcal{R}(V_N) = \mathcal{N}(\mathbf{1}^T), \quad V_N^T \mathbf{1} = \mathbf{0} \quad (1040)$$

Matrix-valued function (1038) meets the criterion for convexity in §3.13.0.0.2 over its domain that is all of  $\mathbb{R}^{n \times N}$ ; *videlicet*, for any  $Y \in \mathbb{R}^{n \times N}$

$$-\frac{d^2}{dt^2} V_N^T \mathbf{D}(X + tY) V_N = 4V_N^T Y^T Y V_N \succeq 0 \quad (1041)$$

Quadratic matrix-valued function  $-V_N^T \mathbf{D}(X) V_N$  is therefore convex in  $X$  achieving its minimum, with respect to a positive semidefinite cone (§2.7.2.2), at  $X = \mathbf{0}$ . When the penultimate number of points exceeds the dimension of the space  $n < N-1$ , strict convexity of the quadratic (1038) becomes impossible because (1041) could not then be positive definite.

### 5.4.2 Gram-form EDM definition

Positive semidefinite matrix  $X^T X$  in (1033), formed from inner product of list  $X$ , is known as a *Gram matrix*; [280, §3.6]

$$\begin{aligned} G \triangleq X^T X &= \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \begin{bmatrix} x_1 & \cdots & x_N \end{bmatrix} = \begin{bmatrix} \|x_1\|^2 & x_1^T x_2 & x_1^T x_3 & \cdots & x_1^T x_N \\ x_2^T x_1 & \|x_2\|^2 & x_2^T x_3 & \cdots & x_2^T x_N \\ x_3^T x_1 & x_3^T x_2 & \|x_3\|^2 & \ddots & x_3^T x_N \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ x_N^T x_1 & x_N^T x_2 & x_N^T x_3 & \cdots & \|x_N\|^2 \end{bmatrix} \in \mathbb{S}_+^N \\ &= \delta \left( \begin{bmatrix} \|x_1\| \\ \|x_2\| \\ \vdots \\ \|x_N\| \end{bmatrix} \right) \begin{bmatrix} 1 & \cos \psi_{12} & \cos \psi_{13} & \cdots & \cos \psi_{1N} \\ \cos \psi_{12} & 1 & \cos \psi_{23} & \cdots & \cos \psi_{2N} \\ \cos \psi_{13} & \cos \psi_{23} & 1 & \ddots & \cos \psi_{3N} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \cos \psi_{1N} & \cos \psi_{2N} & \cos \psi_{3N} & \cdots & 1 \end{bmatrix} \delta \left( \begin{bmatrix} \|x_1\| \\ \|x_2\| \\ \vdots \\ \|x_N\| \end{bmatrix} \right) \quad (1042) \\ &\triangleq \sqrt{\delta^2(G)} \Psi \sqrt{\delta^2(G)} \end{aligned}$$

where  $\psi_{ij}$  (1061) is angle between vectors  $x_i$  and  $x_j$ , and where  $\delta^2$  denotes a diagonal matrix in this case. Positive semidefiniteness of *interpoint angle matrix*  $\Psi$  implies positive semidefiniteness of Gram matrix  $G$ ;

$$G \succeq 0 \Leftrightarrow \Psi \succeq 0 \quad (1043)$$

When  $\delta^2(G)$  is nonsingular, then  $G \succeq 0 \Leftrightarrow \Psi \succeq 0$ . (§A.3.1.0.5)

Distance-square  $d_{ij}$  (1029) is related to Gram matrix entries  $G^T = G \triangleq [g_{ij}]$

$$\begin{aligned} d_{ij} &= g_{ii} + g_{jj} - 2g_{ij} \\ &= \langle \Phi_{ij}, G \rangle \end{aligned} \quad (1044)$$

where  $\Phi_{ij}$  is defined in (1031). Hence the linear EDM definition

$$\left. \begin{aligned} \mathbf{D}(G) &\triangleq \delta(G)\mathbf{1}^T + \mathbf{1}\delta(G)^T - 2G \in \mathbb{EDM}^N \\ &= [\langle \Phi_{ij}, G \rangle, i, j=1 \dots N] \end{aligned} \right\} \Leftrightarrow G \succeq 0 \quad (1045)$$

The EDM cone may be described, (*confer* (1134)(1140))

$$\mathbb{EDM}^N = \left\{ \mathbf{D}(G) \mid G \in \mathbb{S}_+^N \right\} \quad (1046)$$

#### 5.4.2.1 First point at origin

Assume the first point  $x_1$  in an unknown list  $X$  resides at the origin;

$$Xe_1 = \mathbf{0} \Leftrightarrow Ge_1 = \mathbf{0} \quad (1047)$$

Consider the symmetric translation  $(I - \mathbf{1}e_1^T)\mathbf{D}(G)(I - e_1\mathbf{1}^T)$  that shifts the first row and column of  $\mathbf{D}(G)$  to the origin; setting Gram-form EDM operator  $\mathbf{D}(G) = D$  for convenience,

$$-(D - (De_1\mathbf{1}^T + \mathbf{1}e_1^T D) + \mathbf{1}e_1^T De_1\mathbf{1}^T) \frac{1}{2} = G - (Ge_1\mathbf{1}^T + \mathbf{1}e_1^T G) + \mathbf{1}e_1^T Ge_1\mathbf{1}^T \quad (1048)$$

where

$$e_1 \triangleq \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (1049)$$

is the first vector from the standard basis. Then it follows, for  $D \in \mathbb{S}_h^N$

$$\begin{aligned} G &= -(D - (De_1\mathbf{1}^T + \mathbf{1}e_1^T D))^{\frac{1}{2}}, \quad x_1 = \mathbf{0} \\ &= -[\mathbf{0} \ \sqrt{2}V_N]^T D [\mathbf{0} \ \sqrt{2}V_N]^{\frac{1}{2}} \\ &= \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & -V_N^T D V_N \end{bmatrix} \quad (1050) \\ V_N^T G V_N &= -V_N^T D V_N^{\frac{1}{2}} \quad \forall X \end{aligned}$$

where

$$I - e_1\mathbf{1}^T = \begin{bmatrix} \mathbf{0} & \sqrt{2}V_N \end{bmatrix} \quad (1051)$$

is a projector (§B.4.2 no.7) nonorthogonally projecting (§E.1, §E.8) on subspace

$$\begin{aligned} \mathbb{S}_{\mathbf{0}}^N &= \{G \in \mathbb{S}^N \mid Ge_1 = \mathbf{0}\} \\ &= \left\{ [\mathbf{0} \ \sqrt{2}V_N]^T Y [\mathbf{0} \ \sqrt{2}V_N] \mid Y \in \mathbb{S}^N \right\} \quad (2200) \end{aligned}$$

in the Euclidean sense. From (1050) we get sufficiency of the first matrix criterion for an EDM proved by Schoenberg in 1935; [349]<sup>5.7</sup>

$$D \in \text{EDM}^N \Leftrightarrow \begin{cases} -V_N^T D V_N \in \mathbb{S}_+^{N-1} \\ D \in \mathbb{S}_h^N \end{cases} \quad (1052)$$

We provide a rigorous complete more geometric proof of this fundamental *Schoenberg criterion* in §5.9.1.0.4. [431, isedm()]

By substituting  $G = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & -V_N^T D V_N \end{bmatrix}$  (1050) into  $\mathbf{D}(G)$  (1045),

$$D = \begin{bmatrix} 0 \\ \delta(-V_N^T D V_N) \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & \delta(-V_N^T D V_N)^T \end{bmatrix} - 2 \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & -V_N^T D V_N \end{bmatrix} \quad (1154)$$

assuming  $x_1 = \mathbf{0}$ . Details of this bijection are provided in §5.6.2.

#### 5.4.2.2 0 geometric center

Assume the *geometric center* (§5.5.1.0.1) of an unknown list  $X$  is the origin;

$$X\mathbf{1} = \mathbf{0} \Leftrightarrow G\mathbf{1} = \mathbf{0} \quad (1053)$$

Now consider the calculation  $(I - \frac{1}{N}\mathbf{1}\mathbf{1}^T)\mathbf{D}(G)(I - \frac{1}{N}\mathbf{1}\mathbf{1}^T)$ , a geometric centering or projection operation. (§E.7.2.0.2) Setting  $\mathbf{D}(G) = D$  for convenience as in §5.4.2.1,

$$\begin{aligned} G &= -(D - \frac{1}{N}(D\mathbf{1}\mathbf{1}^T + \mathbf{1}\mathbf{1}^T D) + \frac{1}{N^2}\mathbf{1}\mathbf{1}^T D \mathbf{1}\mathbf{1}^T)^{\frac{1}{2}}, \quad X\mathbf{1} = \mathbf{0} \\ &= -VDV^{\frac{1}{2}} \\ VGV &= -VDV^{\frac{1}{2}} \quad \forall X \end{aligned} \quad (1054)$$

---

<sup>5.7</sup>From (1040) we know  $\mathcal{R}(V_N) = \mathcal{N}(\mathbf{1}^T)$ , so (1052) is the same as (1028). In fact, any matrix  $V$  in place of  $V_N$  will satisfy (1052) whenever  $\mathcal{R}(V) = \mathcal{R}(V_N) = \mathcal{N}(\mathbf{1}^T)$ . But  $V_N$  is the matrix implicit in Schoenberg's seminal exposition.

where more properties of the auxiliary (*geometric centering*, projection) matrix

$$V \triangleq I - \frac{1}{N} \mathbf{1}\mathbf{1}^T \in \mathbb{S}^N \quad (1055)$$

are found in §B.4. From (1054) and the assumption  $D \in \mathbb{S}_h^N$  we get sufficiency of the more popular form of Schoenberg criterion:

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} -VDV \in \mathbb{S}_+^N \\ D \in \mathbb{S}_h^N \end{cases} \quad (1056)$$

Of particular utility when  $D \in \mathbb{EDM}^N$  is the fact, (§B.4.2 no.20) (1029)

$$\begin{aligned} \text{tr}(-VDV \frac{1}{2}) &= \frac{1}{2N} \sum_{i,j} d_{ij} &= \frac{1}{2N} \text{vec}(X)^T \left( \sum_{i,j} \Phi_{ij} \otimes I \right) \text{vec } X \\ &= \text{tr}(VGV), \quad G \succeq 0 && (1057) \\ &= \text{tr } G &= \sum_{\ell=1}^N \|x_\ell\|^2 = \|X\|_F^2, \quad X\mathbf{1} = \mathbf{0} \end{aligned}$$

where  $\sum \Phi_{ij} \in \mathbb{S}_+^N$  (1031), therefore convex in  $\text{vec } X$ . We will find this trace useful as a heuristic to minimize affine dimension of an unknown list arranged columnar in  $X$  (§7.2.2), but it tends to facilitate reconstruction of a list configuration having least energy; *id est*, it compacts a reconstructed list by minimizing total norm-square of the vertices.

By substituting  $G = -VDV \frac{1}{2}$  (1054) into  $\mathbf{D}(G)$  (1045), assuming  $X\mathbf{1} = \mathbf{0}$

$$D = \delta(-VDV \frac{1}{2})\mathbf{1}^T + \mathbf{1}\delta(-VDV \frac{1}{2})^T - 2(-VDV \frac{1}{2}) \quad (1144)$$

Details of this bijection can be found in §5.6.1.1.

#### 5.4.2.2.1 Example. Hypersphere.

These foregoing relationships allow combination of distance and Gram constraints in any optimization problem we might pose:

- Interpoint angle  $\Psi$  can be constrained by fixing all individual point lengths  $\sqrt{\delta(G)}$ ; then

$$\Psi = -\sqrt{\delta^2(G)}^{-1} VDV \frac{1}{2} \sqrt{\delta^2(G)}^{-1}, \quad X\mathbf{1} = \mathbf{0} \quad (1058)$$

- (confer §5.9.1.0.3, (1243) (1387)) Constraining all main diagonal entries  $g_{ii}$  of a Gram matrix to 1, for example, is equivalent to the constraint that all points lie on a hypersphere of radius 1 centered at the origin.

$$D = 2(g_{11}\mathbf{1}\mathbf{1}^T - G) \in \mathbb{EDM}^N \quad (1059)$$

Requiring  $\mathbf{0}$  geometric center then becomes equivalent to the constraint:  $D\mathbf{1} = 2N\mathbf{1}$ . [97, p.116] Any further constraint on that Gram matrix applies only to interpoint angle matrix  $\Psi = G$ .

Because any point list may be constrained to lie on a hypersphere boundary whose affine dimension exceeds that of the list, a Gram matrix may always be constrained to have equal positive values along its main diagonal. (Laura Klanfer 1933 [349, §3]) This observation renewed interest in the ellotope (§5.9.1.0.1).  $\square$

#### 5.4.2.2.2 Example. List-member constraints via Gram matrix.

Capitalizing on identity (1054) relating Gram and EDM  $D$  matrices, a constraint set such as

$$\left. \begin{aligned} \text{tr}\left(-\frac{1}{2}V D V e_i e_i^T\right) &= \|x_i\|^2 \\ \text{tr}\left(-\frac{1}{2}V D V (e_i e_j^T + e_j e_i^T)\frac{1}{2}\right) &= x_i^T x_j \\ \text{tr}\left(-\frac{1}{2}V D V e_j e_j^T\right) &= \|x_j\|^2 \end{aligned} \right\} \quad (1060)$$

relates list member  $x_i$  to  $x_j$  to within an isometry through inner-product identity [442, §1-7]

$$\cos \psi_{ij} = \frac{x_i^T x_j}{\|x_i\| \|x_j\|} \quad (1061)$$

where  $\psi_{ij}$  is angle between the two vectors as in (1042). For  $M$  list members, there total  $M(M+1)/2$  such constraints. Angle constraints are incorporated in Example 5.4.2.2.5 and Example 5.4.2.2.13.  $\square$

#### 5.4.2.2.3 Example. Gram matrix as optimization problem.

Consider the academic problem of finding a Gram matrix (1054) subject to constraints on each and every entry of the corresponding EDM:

$$\begin{aligned} &\underset{D \in \mathbb{S}_h^N}{\text{find}} \quad -V D V \frac{1}{2} \in \mathbb{S}^N \\ &\text{subject to} \quad \langle D, (e_i e_j^T + e_j e_i^T)\frac{1}{2} \rangle = \check{d}_{ij}, \quad i, j = 1 \dots N, \quad i < j \\ &\quad -V D V \succeq 0 \end{aligned} \quad (1062)$$

where the  $\check{d}_{ij}$  are given nonnegative constants. EDM  $D$  can, of course, be replaced with the equivalent Gram-form (1045). Requiring only the selfadjointness property (1572) of the main-diagonal linear operator  $\delta$  we get, for  $A \in \mathbb{S}^N$

$$\langle D, A \rangle = \langle \delta(G) \mathbf{1}^T + \mathbf{1} \delta(G)^T - 2G, A \rangle = 2 \langle G, \delta(A \mathbf{1}) - A \rangle \quad (1063)$$

Then the problem equivalent to (1062) becomes, for  $G \in \mathbb{S}_c^N \Leftrightarrow G \mathbf{1} = \mathbf{0}$

$$\begin{aligned} &\underset{G \in \mathbb{S}_c^N}{\text{find}} \quad G \in \mathbb{S}^N \\ &\text{subject to} \quad \left\langle G, \delta((e_i e_j^T + e_j e_i^T) \mathbf{1}) - (e_i e_j^T + e_j e_i^T) \right\rangle = \check{d}_{ij}, \quad i, j = 1 \dots N, \quad i < j \\ &\quad G \succeq 0 \end{aligned} \quad (1064)$$

Barvinok's Proposition 2.9.3.0.1 predicts existence for either formulation (1062) or (1064) such that implicit equality constraints induced by subspace membership are ignored

$$\text{rank } G, \text{ rank } V D V \leq \left\lfloor \frac{\sqrt{8(N(N-1)/2) + 1} - 1}{2} \right\rfloor = N - 1 \quad (1065)$$

because, in each case, the Gram matrix is confined to a face of positive semidefinite cone  $\mathbb{S}_+^N$  isomorphic with  $\mathbb{S}_+^{N-1}$  (§6.6.1). (§E.7.2.0.2) This bound is tight (§5.7.1.1) and is the greatest upper bound.  $\square$

#### 5.4.2.2.4 Example. First duality.

Kuhn reports that the first dual optimization problem<sup>5.9</sup> to be recorded in the literature

<sup>5.8</sup>  $-V D V|_{N \leftarrow 1} = 0$  (§B.4.1)

<sup>5.9</sup> By *dual problem* is meant, in the strongest sense: the optimal objective achieved by a maximization problem, dual to a given (primal) minimization problem, is always equal to the optimal objective achieved by the minimization. (Figure 64 Example 2.13.1.1) A dual problem is always convex when derived from a primal via Lagrangian function.



Figure 143: Rendering of *Fermat point* in acrylic on canvas by [Suman Vaze](#). Three circles intersect at Fermat point of minimum total distance from three vertices of (and interior to) red/black/white triangle.

dates back to 1755. [419] Perhaps more intriguing is the fact: this earliest instance of duality is a two-dimensional Euclidean distance geometry problem known as *Fermat point* (Figure 143) named after the French mathematician. Given  $N$  distinct points in the plane  $\{x_i \in \mathbb{R}^2, i=1 \dots N\}$ , the Fermat point  $y$  is an optimal solution to

$$\underset{y}{\text{minimize}} \quad \sum_{i=1}^N \|y - x_i\| \quad (1066)$$

a convex minimization of total distance. The historically first dual problem formulation asks for the smallest equilateral triangle encompassing ( $N=3$ ) three points  $x_i$ . Another problem dual to (1066) (Kuhn 1967)

$$\begin{aligned} & \underset{\{z_i\}}{\text{maximize}} \quad \sum_{i=1}^N \langle z_i, x_i \rangle \\ & \text{subject to} \quad \sum_{i=1}^N z_i = \mathbf{0} \\ & \quad \|z_i\| \leq 1 \quad \forall i \end{aligned} \quad (1067)$$

has interpretation as minimization of work required to balance potential energy in an  $N$ -way tug-of-war between equally matched opponents situated at  $\{x_i\}$ . [437]

It is not so straightforward to write the Fermat point problem (1066) equivalently in terms of a Gram matrix from this section. Squaring instead

$$\underset{\alpha}{\text{minimize}} \sum_{i=1}^N \|\alpha - x_i\|^2 \equiv \underset{D \in \mathbb{S}^{N+1}}{\text{minimize}} \quad \langle -V, D \rangle \quad \text{subject to} \quad \begin{aligned} & \langle D, e_i e_j^T + e_j e_i^T \rangle \frac{1}{2} = \check{d}_{ij} \quad \forall (i, j) \in \mathcal{I} \\ & -V D V \succeq 0 \end{aligned} \quad (1068)$$

yields an inequivalent convex geometric centering problem whose equality constraints comprise EDM  $D$  main-diagonal zeros and known distances-square.<sup>5.10</sup> Going the other

<sup>5.10</sup>  $\alpha^*$  is geometric center of points  $x_i$  (1118). For three points,  $\mathcal{I} = \{1, 2, 3\}$ ; optimal affine dimension

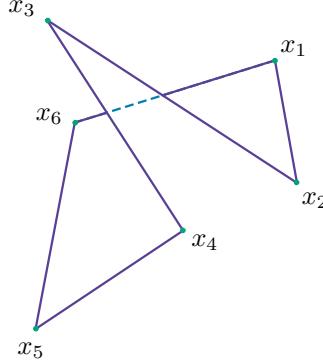


Figure 144: Arbitrary hexagon in  $\mathbb{R}^3$  whose vertices are labelled clockwise.

way, a problem dual to total distance-square maximization (Example 6.7.0.0.1) is a penultimate minimum eigenvalue problem having application to *PageRank* calculation by search engines [259, §4]. [377]

Fermat function (1066) is empirically compared with (1068) in [65, §8.7.3], but for multiple unknowns in  $\mathbb{R}^2$ , where propensity of (1066) for producing zero distance between unknowns is revealed. An optimal solution to (1066) gravitates toward gradient discontinuities (§D.2.1), as in Figure 76, whereas optimal solution to (1068) is less compact in the unknowns. <sup>5.11</sup>  $\square$

#### 5.4.2.2.5 Example. Hexagon.

Barvinok [26, §2.6] poses a problem in *geometric realizability* of an arbitrary hexagon (Figure 144) having:

1. prescribed (one-dimensional) face-lengths  $l$
2. prescribed angles  $\varphi$  between the three pairs of opposing faces
3. a constraint on the sum of norm-square of each and every vertex  $x$ ;

ten affine equality constraints in all on a Gram matrix  $G \in \mathbb{S}^6$  (1054). Let's realize this as a convex feasibility problem (with constraints written in the same order) also assuming  $\mathbf{0}$  geometric center (1053):

$$\begin{aligned} & \underset{D \in \mathbb{S}_h^6}{\text{find}} \quad -V D V \frac{1}{2} \in \mathbb{S}^6 \\ & \text{subject to} \quad \text{tr}(D(e_i e_j^T + e_j e_i^T) \frac{1}{2}) = l_{ij}^2, \quad j-1 = (i=1 \dots 6) \bmod 6 \\ & \quad \text{tr}(-\frac{1}{2} V D V (A_i + A_i^T) \frac{1}{2}) = \cos \varphi_i, \quad i=1, 2, 3 \\ & \quad \text{tr}(-\frac{1}{2} V D V) = 1 \\ & \quad -V D V \succeq 0 \end{aligned} \tag{1069}$$

where, for  $A_i \in \mathbb{R}^{6 \times 6}$  (1061)

$$\begin{aligned} A_1 &= (e_1 - e_6)(e_3 - e_4)^T / (l_{61} l_{34}) \\ A_2 &= (e_2 - e_1)(e_4 - e_5)^T / (l_{12} l_{45}) \\ A_3 &= (e_3 - e_2)(e_5 - e_6)^T / (l_{23} l_{56}) \end{aligned} \tag{1070}$$

(§5.7) must be 2 because a third dimension can only increase total distance. Minimization of  $\langle -V, D \rangle$  is a heuristic for rank minimization. (§7.2.2)

<sup>5.11</sup>Optimal solution to (1066) has mechanical interpretation in terms of interconnecting springs with constant force when distance is nonzero; otherwise, 0 force. Problem (1068) is interpreted instead using linear springs.

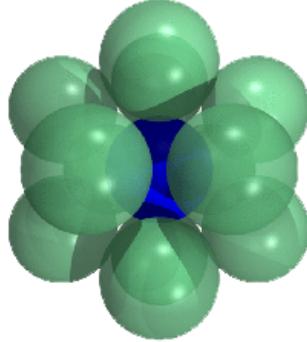


Figure 145: Sphere-packing illustration from [436, *kissing number*]. Translucent balls illustrated all have the same diameter.

and where the first constraint on length-square  $l_{ij}^2$  can be equivalently written as a constraint on the Gram matrix  $-VDV\frac{1}{2}$  via (1063). We show how to numerically solve such a problem by *alternating projection* in §E.10.2.1.1. Barvinok's Proposition 2.9.3.0.1 asserts existence of a list, corresponding to Gram matrix  $G$  solving this feasibility problem, whose affine dimension (§5.7.1.1) does not exceed 3 because the convex feasible set is bounded by the third constraint  $\text{tr}(-\frac{1}{2}VDV)=1$  (1057).  $\square$

#### 5.4.2.2.6 Example. Kissing number of sphere packing.

Two nonoverlapping Euclidean balls are said to *kiss* if they touch. An elementary geometrical problem can be posed: *Given hyperspheres, each having the same diameter 1, how many hyperspheres can simultaneously kiss one central hypersphere?* [457] Noncentral hyperspheres are allowed, but not required, to kiss.

As posed, the problem seeks the maximal number of spheres  $K$  kissing a central sphere in a particular dimension. The total number of spheres is  $N = K + 1$ . In one dimension, the answer to this kissing problem is 2. In two dimensions, 6. (Figure 9)

The question was presented, in three dimensions, to Isaac Newton by David Gregory in the context of celestial mechanics. And so was born a controversy between the two scholars on the campus of Trinity College Cambridge in 1694. Newton correctly identified the kissing number as 12 (Figure 145) while Gregory argued for 13. Their dispute was finally resolved in 1953 by Schütte & van der Waerden. [334] In 2003, Oleg Musin tightened the upper bound on kissing number  $K$  in four dimensions from 25 to  $K = 24$  by refining a method of Philippe Delsarte from 1973. Delsarte's method provides an infinite number [17] of linear inequalities necessary for converting a rank-constrained semidefinite program<sup>5.12</sup> to a linear program.<sup>5.13</sup> [304]

There are no proofs known for kissing number in higher dimension excepting dimensions eight and twenty four. Interest persists [92] because sphere packing has found application to error correcting codes from the fields of communications and information theory; specifically to quantum computing. [100]

Translating this problem to an *EDM graph* realization (Figure 141, Figure 146) is suggested by Pfender & Ziegler. Imagine the centers of each sphere are connected by line

<sup>5.12</sup>whose feasible set belongs to that subset of an ellotope (§5.9.1.0.1) bounded above by some desired rank.

<sup>5.13</sup>Simplex-method solvers for linear programs produce numerically better results than contemporary log-barrier (interior-point method) solvers, for semidefinite programs, by about 7 orders of magnitude; they are far more predisposed to vertex solutions [103, p.158].

segments. Then the distance between centers must obey simple criteria: Each sphere touching the central sphere has a line segment of length exactly 1 joining its center to the central sphere's center. All spheres, excepting the central sphere, must have centers separated by a distance of at least 1.

From this perspective, the kissing problem can be posed as a semidefinite program. Assign index 1 to the central sphere assuming a total of  $N$  spheres:

$$\begin{aligned} & \underset{D \in \mathbb{S}^N}{\text{minimize}} && -\text{tr}(WV_N^T DV_N) \\ & \text{subject to} && D_{1j} = 1, \quad j = 2 \dots N \\ & && D_{ij} \geq 1, \quad 2 \leq i < j = 3 \dots N \\ & && D \in \mathbb{EDM}^N \end{aligned} \quad (1071)$$

Then kissing number

$$K = N_{\max} - 1 \quad (1072)$$

is found from the maximal number  $N$  of spheres that solve this semidefinite program in a given affine dimension  $r$  whose realization is assured by 0 optimal objective. Matrix  $W$  is constant, in this program, determined by a method disclosed in §4.5.1. Matrix  $W \in \mathbb{S}_+^{N-1}$  can be interpreted as direction of search through the positive semidefinite cone for a rank- $r$  optimal solution  $-V_N^T D^* V_N \in \mathbb{S}_+^{N-1}$ : In one dimension, optimal direction matrix  $W^*$  has rank  $= K - r = 2 - 1 = 1$ ;

$$W^* = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \frac{1}{2} \quad (1073)$$

In two dimensions, optimal  $W^*$  has rank  $= K - r = 6 - 2 = 4$ ;

$$W^* = \begin{bmatrix} 4 & 1 & 2 & -1 & -1 & 1 \\ 1 & 4 & -1 & -1 & 2 & 1 \\ 2 & -1 & 4 & 1 & 1 & -1 \\ -1 & -1 & 1 & 4 & 1 & 2 \\ -1 & 2 & 1 & 1 & 4 & -1 \\ 1 & 1 & -1 & 2 & -1 & 4 \end{bmatrix} \frac{1}{6} \quad (1074)$$

In three dimensions, we leave it an exercise to find a rational optimal direction matrix  $W^*$  having rank  $= K - r = 12 - 3 = 9$ . Here is a full-rank rational optimal direction matrix:

$$W^* = \begin{bmatrix} 9 & 1 & -2 & -1 & 3 & -1 & -1 & 1 & 2 & 1 & -2 & 1 \\ 1 & 9 & 3 & -1 & -1 & 1 & 1 & -2 & 1 & 2 & -1 & -1 \\ -2 & 3 & 9 & 1 & 2 & -1 & -1 & 2 & -1 & -1 & 1 & 2 \\ -1 & -1 & 1 & 9 & 1 & -1 & 1 & -1 & 3 & 2 & -1 & 1 \\ 3 & -1 & 2 & 1 & 9 & 1 & 1 & -1 & -1 & -1 & 1 & -1 \\ -1 & 1 & -1 & -1 & 1 & 9 & 2 & -1 & 2 & -1 & 2 & 3 \\ -1 & 1 & -1 & 1 & 1 & 2 & 9 & 3 & -1 & 1 & -2 & -1 \\ 1 & -2 & 2 & -1 & -1 & -1 & 3 & 9 & 2 & -1 & 1 & 1 \\ 2 & 1 & -1 & 3 & -1 & 2 & -1 & 2 & 9 & -1 & 1 & -1 \\ 1 & 2 & -1 & 2 & -1 & -1 & 1 & -1 & -1 & 9 & 3 & 1 \\ -2 & -1 & 1 & -1 & 1 & 2 & -2 & 1 & 1 & 3 & 9 & -1 \\ 1 & -1 & 2 & 1 & -1 & 3 & -1 & 1 & -1 & 1 & -1 & 9 \end{bmatrix} \frac{1}{12} \quad (1075)$$

A four-dimensional solution also has rational optimal direction matrix  $W^*$  having rank  $= K - r = 24 - 4 = 20$ ;

$$W^* = \begin{bmatrix} 20 & -2 & 2 & -2 & 0 & 0 & -2 & 2 & 2 & -2 & 2 & 0 & 2 & 4 & -2 & 2 & 0 & -2 & -2 & 2 & 0 & 0 & 0 & -2 & 2 & -2 \\ -2 & 20 & 2 & 0 & 2 & -2 & -2 & 0 & 2 & 0 & -2 & 2 & -2 & -2 & 0 & -2 & -2 & 0 & 0 & 2 & -2 & -2 & 0 & 0 & -2 & 2 & -2 \\ 2 & 2 & 20 & 2 & 2 & 2 & 0 & 0 & 2 & 0 & -2 & 0 & 2 & -2 & 0 & -2 & -2 & 0 & 0 & 2 & -2 & -2 & 0 & 0 & 2 & 2 & -2 \\ -2 & 0 & 2 & 20 & -2 & 2 & -2 & 0 & -2 & 0 & 2 & -2 & 4 & 2 & 2 & 0 & -2 & 2 & 0 & -2 & 0 & 0 & 2 & 2 & 2 & -2 \\ 0 & 2 & 2 & -2 & 20 & 0 & 2 & -2 & -2 & 2 & -2 & 0 & 2 & 0 & 2 & -2 & 0 & 2 & -2 & 2 & 4 & 0 & 0 & -2 & 2 & -2 \\ 0 & -2 & 2 & 2 & 0 & 20 & 2 & -2 & 2 & 2 & -2 & 0 & -2 & 0 & -2 & 2 & 4 & -2 & 2 & -2 & 2 & 0 & 0 & 0 & -2 & 0 \\ -2 & -2 & 0 & -2 & 2 & 2 & 20 & 2 & 0 & -2 & 4 & -2 & 2 & 2 & 0 & 2 & -2 & 0 & 0 & 2 & -2 & -2 & 2 & 0 & 0 & -2 & 0 \\ -2 & 0 & 2 & 0 & -2 & -2 & 2 & 20 & -2 & 4 & -2 & -2 & 0 & -2 & -2 & 0 & 2 & 2 & 2 & 0 & 0 & 2 & 2 & 2 & -2 & 0 \\ 2 & 2 & 0 & -2 & -2 & 2 & 2 & 0 & 2 & 0 & -2 & 2 & 2 & -2 & 0 & -2 & -2 & 4 & 0 & 2 & 2 & 2 & 2 & 0 & 0 & 2 & 0 \\ -2 & 0 & -2 & 0 & 2 & 2 & -2 & 4 & 2 & 20 & 2 & 2 & 0 & 2 & 2 & 0 & -2 & -2 & 0 & 0 & -2 & -2 & 0 & 0 & -2 & 2 & 0 \\ 2 & 2 & 2 & 0 & -2 & -2 & -2 & 4 & 0 & 2 & 20 & 2 & -2 & 0 & -2 & 2 & 0 & 0 & -2 & 2 & 2 & -2 & 0 & 0 & -2 & 0 & -2 \\ 0 & -2 & 2 & -2 & 0 & 0 & -2 & -2 & -2 & 2 & 20 & 2 & 0 & -2 & 2 & 0 & 2 & 2 & 2 & -2 & 0 & 0 & 4 & -2 & -2 & 0 \\ 2 & 0 & -2 & 4 & 2 & -2 & 2 & 0 & 0 & -2 & 2 & 20 & -2 & -2 & 0 & 2 & -2 & 2 & 0 & 0 & -2 & -2 & 2 & 0 & 0 & -2 & 2 \\ 4 & 2 & -2 & 2 & 0 & 0 & 2 & -2 & -2 & 2 & -2 & 0 & -2 & 20 & -2 & -2 & 0 & 2 & 2 & 2 & -2 & 0 & 0 & 0 & 2 & 0 \\ 2 & -2 & 0 & 2 & 2 & -2 & 0 & -2 & 0 & 2 & 0 & -2 & -2 & 20 & 2 & 2 & 0 & 4 & 2 & -2 & -2 & 2 & 0 & 0 & 2 & 0 \\ 2 & 4 & -2 & 0 & -2 & 2 & 2 & 0 & -2 & 0 & -2 & 2 & 0 & -2 & 20 & 2 & -2 & 0 & 0 & 2 & -2 & 0 & 0 & 2 & -2 & 2 & 0 \\ 0 & 2 & -2 & -2 & 0 & 4 & -2 & -2 & 2 & -2 & 2 & 0 & 2 & 0 & 2 & -2 & 20 & 2 & -2 & 0 & 0 & 0 & 0 & 0 & 0 & 2 \\ -2 & -2 & 0 & 2 & 2 & -2 & 0 & 2 & 4 & -2 & 0 & -2 & 0 & 2 & -2 & 2 & 0 & 2 & 20 & 0 & -2 & 2 & -2 & -2 & 0 & 0 \\ -2 & 2 & 0 & -2 & -2 & 2 & 0 & 2 & 0 & -2 & 0 & 2 & 2 & 4 & -2 & -2 & 0 & 0 & 20 & -2 & 2 & 2 & -2 & 0 & 0 \\ -2 & 0 & 2 & 0 & -2 & -2 & 2 & 0 & 0 & -2 & 2 & 0 & 0 & 2 & -2 & -2 & 0 & 2 & -2 & 20 & 4 & 2 & -2 & -2 & 0 \\ 0 & -2 & -2 & 2 & 4 & 0 & -2 & 2 & 2 & -2 & 2 & 0 & -2 & 0 & -2 & 2 & 0 & 0 & -2 & 2 & 2 & -2 & 20 & 0 & 0 \\ 0 & 2 & -2 & 2 & 0 & 0 & 2 & 2 & -2 & -2 & 4 & -2 & 0 & 2 & -2 & 2 & 0 & 0 & 20 & 0 & 20 & 0 & 20 & 0 \\ -2 & -2 & 4 & -2 & -2 & -2 & 0 & 0 & 2 & 0 & -2 & 2 & 2 & 0 & 0 & 2 & 0 & 0 & -2 & 2 & 2 & 0 & 20 & 2 & 20 \end{bmatrix} \frac{1}{24}$$

but these direction matrices are not unique and their precision not critical. Here is an optimal four-dimensional point list, 5.14 in MATLAB output format, reconstructed by a method in §5.12:

<b>Columns 1 through 6</b>						
X =	0	-0.1983	-0.4584	0.1657	0.9399	0.7416
	0	0.6863	0.2936	0.6239	-0.2936	0.3927
	0	-0.4835	0.8146	-0.6448	0.0611	-0.4224
	0	0.5059	0.2004	-0.4093	-0.1632	0.3427
<b>Columns 7 through 12</b>						
-0.4815	-0.9399	-0.7416	0.1983	0.4584	-0.2832	
0	0.2936	-0.3927	-0.6863	-0.2936	-0.6863	
-0.8756	-0.0611	0.4224	0.4835	-0.8146	-0.3922	
-0.0372	0.1632	-0.3427	-0.5059	-0.2004	-0.5431	
<b>Columns 13 through 18</b>						
0.2832	-0.2926	-0.6473	0.0943	0.3640	-0.3640	
0.6863	0.9176	-0.6239	-0.2313	-0.0624	0.0624	
0.3922	0.1698	-0.2309	-0.6533	-0.1613	0.1613	
0.5431	-0.2088	0.3721	0.7147	-0.9152	0.9152	
<b>Columns 19 through 25</b>						
-0.0943	0.6473	-0.1657	0.2926	-0.5759	0.5759	0.4815
0.2313	0.6239	-0.6239	-0.9176	0.2313	-0.2313	0
0.6533	0.2309	0.6448	-0.1698	-0.2224	0.2224	0.8756
-0.7147	-0.3721	0.4093	0.2088	-0.7520	0.7520	0.0372

The  $r$  nonzero optimal eigenvalues of  $-V_N^T D^* V_N$  are equal; remaining eigenvalues are zero as per  $-\text{tr}(W^* V_N^T D^* V_N) = 0$  (802). Numerical problems begin to arise with matrices of this size due to interior-point methods of solution to (1071). By eliminating some equality constraints from the kissing number problem, matrix size can be reduced: From

5.14 An optimal five-dimensional point list is known: *The answer was known at least 175 years ago. I believe Gauss knew it.* Moreover, Korkine & Zolotarev proved in 1882 that  $D_5$  is the densest lattice in five dimensions. So they proved that if a kissing arrangement in five dimensions can be extended to some lattice, then  $k(5)=40$ . Of course, the conjecture in the general case also is:  $k(5)=40$ . You would like to see coordinates? Easily. Let  $A=\sqrt{2}$ . Then  $p(1)=(A, A, 0, 0, 0)$ ,  $p(2)=(-A, A, 0, 0, 0)$ ,  $p(3)=(A, -A, 0, 0, 0)$ , ...  $p(40)=(0, 0, 0, -A, -A)$ ; i.e., we are considering points with coordinates that have two  $A$  and three 0 with any choice of signs and any ordering of the coordinates; the same coordinates-expression in dimensions 3 and 4.

The first miracle happens in dimension 6. There are better packings than  $D_6$  (Conjecture:  $k(6)=72$ ). It's a real miracle how dense the packing is in eight dimensions ( $E_8$ =Korkine & Zolotarev packing that was discovered in 1880s) and especially in dimension 24, that is the so-called Leech lattice.

Actually, people in coding theory have conjectures on the kissing numbers for dimensions up to 32 (or even greater?). However, sometimes they found better lower bounds. I know that Ericson & Zinoviev a few years ago discovered (by hand, no computer) in dimensions 13 and 14 better kissing arrangements than were known before.

—Oleg Musin

§5.8.3 we have

$$-V_N^T D V_N = \mathbf{1} \mathbf{1}^T - [\mathbf{0} \ I] D \begin{bmatrix} \mathbf{0}^T \\ I \end{bmatrix} \frac{1}{2} \quad (1077)$$

(which does not hold more generally) where Identity matrix  $I \in \mathbb{S}^{N-1}$  has one less dimension than EDM  $D$ . By defining an EDM principal submatrix

$$\hat{D} \triangleq [\mathbf{0} \ I] D \begin{bmatrix} \mathbf{0}^T \\ I \end{bmatrix} \in \mathbb{S}_h^{N-1} \quad (1078)$$

so that

$$-V_N^T D V_N = \mathbf{1} \mathbf{1}^T - \hat{D} \frac{1}{2} \quad (1079)$$

we get a convex problem equivalent to (1071)

$$\begin{aligned} & \underset{\hat{D} \in \mathbb{S}^K}{\text{minimize}} && -\text{tr}(W \hat{D}) \\ & \text{subject to} && \hat{D}_{ij} \geq 1, \quad 1 \leq i < j = 2 \dots K \\ & && \mathbf{1} \mathbf{1}^T - \hat{D} \frac{1}{2} \succeq 0 \\ & && \delta(\hat{D}) = \mathbf{0} \end{aligned} \quad (1080)$$

Any feasible solution  $\mathbf{1} \mathbf{1}^T - \hat{D} \frac{1}{2}$  belongs to an ellipope (§5.9.1.0.1).  $\square$

#### 5.4.2.2.7 Exercise. Rational optimal kissing direction matrix $W^*$ .

Replace (1075) with a rational  $W^*$  having rank  $= K - r = 12 - 3 = 9$ , main diagonal 9, and common denominator 12.  $\blacktriangledown$

This next example shows how finding the common point of intersection for three circles in a plane, a nonlinear problem, has convex expression.

#### 5.4.2.2.8 Example. Trilateration in wireless sensor network.

[185]

Given three known absolute point positions in  $\mathbb{R}^2$  (three anchors  $\check{x}_2, \check{x}_3, \check{x}_4$ ) and only one unknown point (one sensor  $x_1$ ), the sensor's absolute position is determined from its noiseless measured distance-square  $\check{d}_{i1}$  to each of three anchors (Figure 4, Figure 146a). This trilateration can be expressed as a convex optimization problem in terms of list  $X \triangleq [x_1 \ \check{x}_2 \ \check{x}_3 \ \check{x}_4] \in \mathbb{R}^{2 \times 4}$  and Gram matrix  $G \in \mathbb{S}^4$  (1042):

$$\begin{aligned} & \underset{G \in \mathbb{S}^4, X \in \mathbb{R}^{2 \times 4}}{\text{minimize}} && \text{tr } G \\ & \text{subject to} && \begin{aligned} \text{tr}(G \Phi_{i1}) &= \check{d}_{i1}, & i &= 2, 3, 4 \\ \text{tr}(G e_i e_i^T) &= \|\check{x}_i\|^2, & i &= 2, 3, 4 \\ \text{tr}(G(e_i e_j^T + e_j e_i^T)/2) &= \check{x}_i^T \check{x}_j, & 2 \leq i < j &= 3, 4 \\ X(:, 2:4) &= [\check{x}_2 \ \check{x}_3 \ \check{x}_4] \\ \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} &\succeq 0 \end{aligned} \end{aligned} \quad (1081)$$

where

$$\Phi_{ij} = (e_i - e_j)(e_i - e_j)^T \in \mathbb{S}_+^n \quad (1031)$$

and where the constraint on distance-square  $\check{d}_{i1}$  is equivalently written as a constraint on the Gram matrix via (1044). There are 9 linearly independent affine equality constraints on that Gram matrix while the sensor is constrained, only by dimensioning, to lie in  $\mathbb{R}^2$ . Although the objective  $\text{tr } G$  of minimization<sup>5.15</sup> insures a solution on the boundary of

<sup>5.15</sup>Trace ( $\text{tr } G = \langle I, G \rangle$ ) minimization is a heuristic for rank minimization. (§7.2.2.1) It may be interpreted as squashing  $G$  which is bounded below by  $X^T X$  as in (1082); *id est*,  $G - X^T X \succeq 0 \Rightarrow \text{tr } G \geq \text{tr } X^T X$  (1643).  $\delta(G - X^T X) = \mathbf{0} \Leftrightarrow G = X^T X$  (§A.7.2)  $\Rightarrow \text{tr } G = \text{tr } X^T X$  which is a condition necessary for equality.

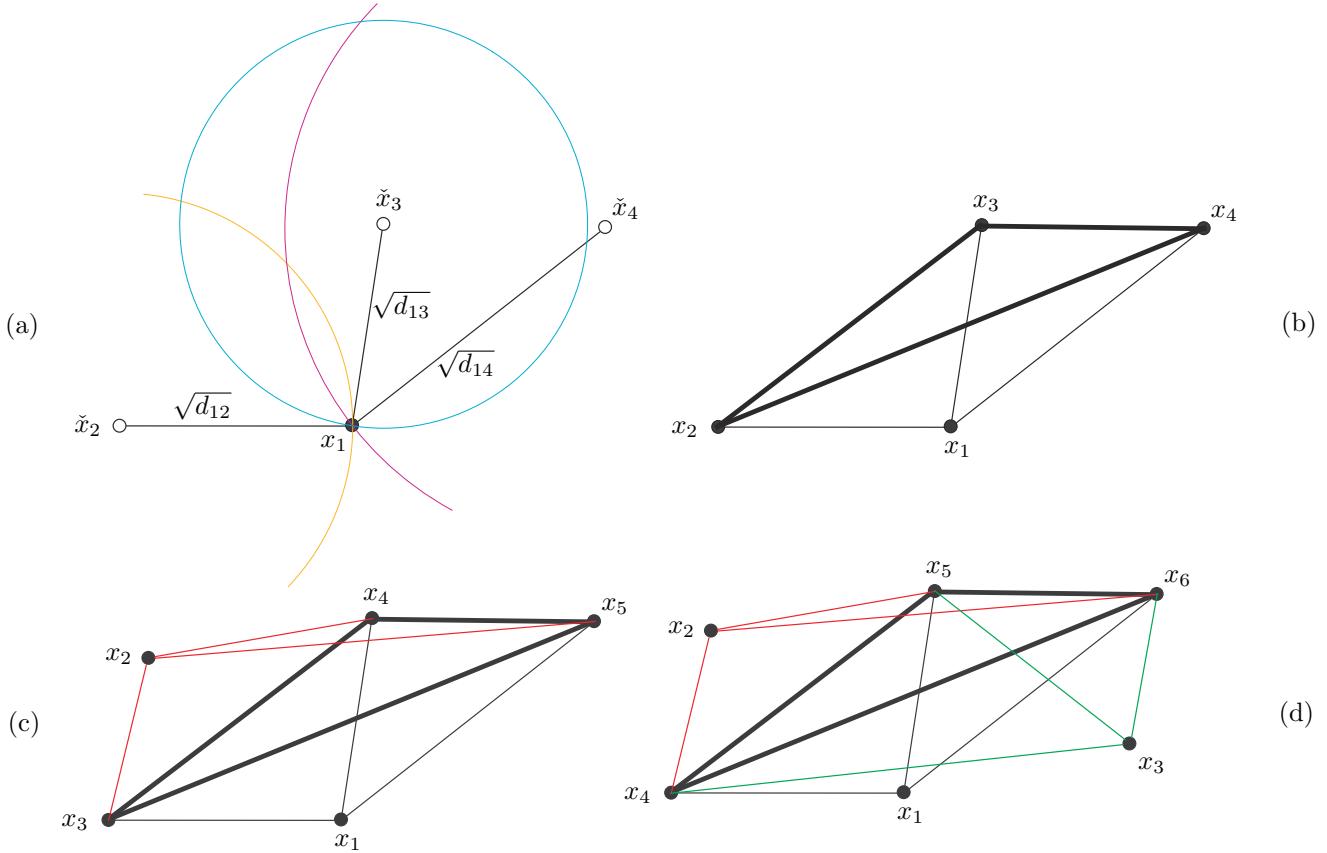


Figure 146: (a) Given three distances indicated with absolute point positions  $\check{x}_2, \check{x}_3, \check{x}_4$  known and noncollinear, absolute position of  $x_1$  in  $\mathbb{R}^2$  can be precisely and uniquely determined by *trilateration*; solution to a system of nonlinear equations. Dimensionless EDM graphs (b) (c) (d) represent EDMs in various states of completion. Line segments represent known absolute distances and may cross without vertex at intersection. (b) Four-point list must always be embeddable in affine subset having dimension  $\text{rank } V_N^T D V_N$  not exceeding 3. To determine relative position of  $x_2, x_3, x_4$ , triangle inequality is necessary and sufficient (§5.14.1). Knowing all distance information, then (by injectivity of  $\mathbf{D}$  (§5.6)) point position  $x_1$  is uniquely determined to within an isometry in any dimension. (c) When fifth point is introduced, only distances to  $x_3, x_4, x_5$  are required to determine relative position of  $x_2$  in  $\mathbb{R}^2$ . Graph represents first instance of missing distance information;  $\sqrt{d_{12}}$ . (d) Three distances are absent ( $\sqrt{d_{12}}, \sqrt{d_{13}}, \sqrt{d_{23}}$ ) from complete set of interpoint distances, yet unique isometric reconstruction (§5.4.2.2.10) of six points in  $\mathbb{R}^2$  is certain.

positive semidefinite cone  $\mathbb{S}_+^4$ , for this problem, we claim that the set of feasible Gram matrices forms a line (§2.5.1.1) in isomorphic  $\mathbb{R}^{10}$  tangent (§2.1.7.1.2) to the positive semidefinite cone boundary. (§5.4.2.2.9, confer §4.2.1.3)

By Schur complement (§A.4, §2.9.1.0.1), any feasible  $G$  and  $X$  provide

$$G \succeq X^T X \quad (1082)$$

which is a convex relaxation of the desired (nonconvex) equality constraint

$$\begin{bmatrix} I & X \\ X^T & G \end{bmatrix} = \begin{bmatrix} I \\ X^T \end{bmatrix} [I \ X] \quad (1083)$$

expected positive semidefinite rank-2 under noiseless conditions. But, by (1645), the relaxation admits

$$(3 \geq) \operatorname{rank} G \geq \operatorname{rank} X \quad (1084)$$

(a third dimension corresponding to an intersection of three spheres, not circles, were there noise). If rank of an optimal solution equals 2,

$$\operatorname{rank} \begin{bmatrix} I & X^* \\ X^{*\top} & G^* \end{bmatrix} = 2 \quad (1085)$$

then  $G^* = X^{*\top} X^*$  by Theorem A.4.0.1.3.

As posed, this *localization* problem does not require affinely independent (Figure 30, three noncollinear) anchors. Assuming the anchors exhibit no rotational or reflective symmetry in their affine hull (§5.5.2) and assuming the sensor  $x_1$  lies in that affine hull, then sensor position solution  $x_1^* = X^*(\cdot, 1)$  is unique under noiseless measurement. [357]

□

This preceding transformation of trilateration to a semidefinite program works all the time ((1085) holds) despite relaxation (1082) because the optimal solution set is a unique point.

**5.4.2.2.9 Proof (sketch).** Only the sensor location  $x_1$  is unknown. The objective function together with the equality constraints make a linear system of equations in Gram matrix variable  $G$

$$\begin{aligned} \operatorname{tr} G &= \|x_1\|^2 + \|\check{x}_2\|^2 + \|\check{x}_3\|^2 + \|\check{x}_4\|^2 \\ \operatorname{tr}(G\Phi_{ii}) &= \check{d}_{i1}, & i &= 2, 3, 4 \\ \operatorname{tr}(G e_i e_i^T) &= \|\check{x}_i\|^2, & i &= 2, 3, 4 \\ \operatorname{tr}(G(e_i e_j^T + e_j e_i^T)/2) &= \check{x}_i^T \check{x}_j, & 2 \leq i < j &= 3, 4 \end{aligned} \quad (1086)$$

which is invertible:

$$\operatorname{svec} G = \left[ \begin{array}{c} \operatorname{svec}(I)^T \\ \operatorname{svec}(\Phi_{21})^T \\ \operatorname{svec}(\Phi_{31})^T \\ \operatorname{svec}(\Phi_{41})^T \\ \operatorname{svec}(e_2 e_2^T)^T \\ \operatorname{svec}(e_3 e_3^T)^T \\ \operatorname{svec}(e_4 e_4^T)^T \\ \operatorname{svec}((e_2 e_3^T + e_3 e_2^T)/2)^T \\ \operatorname{svec}((e_2 e_4^T + e_4 e_2^T)/2)^T \\ \operatorname{svec}((e_3 e_4^T + e_4 e_3^T)/2)^T \end{array} \right]^{-1} \left[ \begin{array}{c} \|x_1\|^2 + \|\check{x}_2\|^2 + \|\check{x}_3\|^2 + \|\check{x}_4\|^2 \\ \check{d}_{21} \\ \check{d}_{31} \\ \check{d}_{41} \\ \|\check{x}_2\|^2 \\ \|\check{x}_3\|^2 \\ \|\check{x}_4\|^2 \\ \check{x}_2^T \check{x}_3 \\ \check{x}_2^T \check{x}_4 \\ \check{x}_3^T \check{x}_4 \end{array} \right] \quad (1087)$$

That line in the ambient space  $\mathbb{S}^4$  of  $G$ , claimed on page 359, is traced by  $\|x_1\|^2 \in \mathbb{R}$  on the right-hand side, as it turns out. One must show this line to be tangential (§2.1.7.1.2) to  $\mathbb{S}_+^4$  in order to prove uniqueness. Tangency is possible for affine dimension 1 or 2 while its occurrence depends completely on the known measurement data. ■

But as soon as significant noise is introduced or whenever distance data is incomplete, such problems can remain convex although the set of all optimal solutions generally becomes a convex set bigger than a single point (and still containing the noiseless solution).

**5.4.2.2.10 Definition.** *Isometric reconstruction.* (confer §5.5.3)

Isometric reconstruction from an EDM means building a list  $X$  correct to within a rotation, reflection, and translation; in other terms, reconstruction of relative position, unique to within an isometry, correct to within a rigid transformation. △

How much distance information is needed to uniquely localize a sensor (to recover actual relative position)? The narrative in Figure 146 helps dispel any notion of distance data proliferation in *low affine dimension* ( $r < N - 2$ ).<sup>5.16</sup> Huang, Liang, and Pardalos [232, §4.2] claim  $O(2N)$  distances is a least lower bound (independent of affine dimension  $r$ ) for unique isometric reconstruction; achievable under certain noiseless conditions on graph connectivity and point position. Alfakih shows how to ascertain uniqueness over all affine dimensions via *Gale matrix*. [10] [5] [6] Figure 141b (p.344, from *small completion problem* Example 5.3.0.0.2) is an example in  $\mathbb{R}^2$  requiring only  $2N - 3 = 5$  known symmetric entries for unique isometric reconstruction, although the four-point example in Figure 146b will not yield a unique reconstruction when any one of the distances is left unspecified.

The list represented by the particular dimensionless *EDM graph* in Figure 147, having only  $2N - 3 = 9$  absolute distances specified, has only one realization in  $\mathbb{R}^2$  but has more realizations in higher dimensions. Unique  $r$ -dimensional isometric reconstruction by semidefinite relaxation like (1082) occurs iff realization in  $\mathbb{R}^r$  is unique and there exist no nontrivial higher-dimensional realizations. [357] For sake of reference, we provide the complete corresponding EDM:

$$D = \begin{bmatrix} 0 & 50641 & 56129 & 8245 & 18457 & 26645 \\ 50641 & 0 & 49300 & 25994 & 8810 & 20612 \\ 56129 & 49300 & 0 & 24202 & 31330 & 9160 \\ 8245 & 25994 & 24202 & 0 & 4680 & 5290 \\ 18457 & 8810 & 31330 & 4680 & 0 & 6658 \\ 26645 & 20612 & 9160 & 5290 & 6658 & 0 \end{bmatrix} \quad (1088)$$

We consider paucity of distance information in this next example which shows it is possible to recover exact relative position given incomplete noiseless distance information. An *ad hoc* method for recovery of the least-rank optimal solution under noiseless conditions is introduced:

---

<sup>5.16</sup>When affine dimension  $r$  reaches  $N - 2$ , then all distances-square in the EDM must be known for unique isometric reconstruction in  $\mathbb{R}^r$ ; going the other way, when  $r = 1$  then the condition that the dimensionless EDM graph be connected is necessary and sufficient. [213, §2.2]

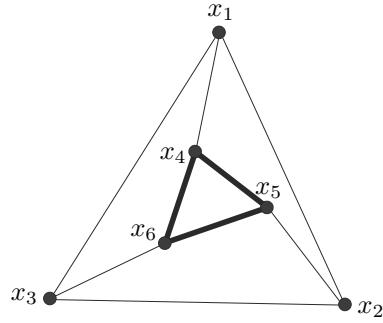


Figure 147: ([confer\(1088\)](#)) Incomplete EDM corresponding to this dimensionless EDM graph (drawn freehand; no symmetry intended) provides unique isometric reconstruction in  $\mathbb{R}^2$ .

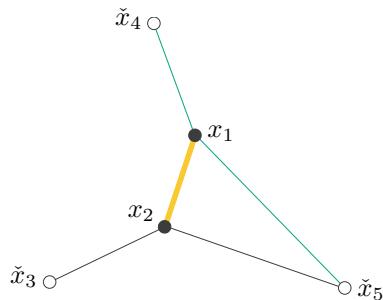


Figure 148: (Ye) Two sensors • and three anchors ○ in  $\mathbb{R}^2$ . Connecting line-segments denote known absolute distances. Incomplete EDM corresponding to this dimensionless EDM graph provides unique isometric reconstruction in  $\mathbb{R}^2$ .

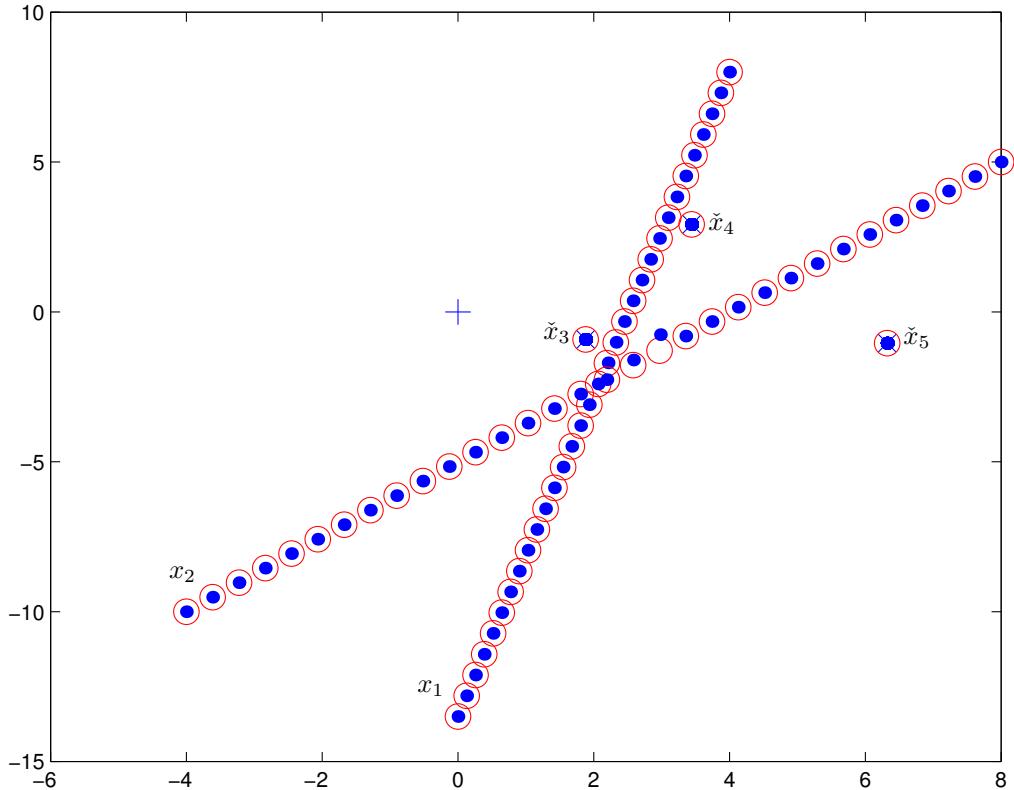


Figure 149: Given in red  $\circ$  are two discrete linear trajectories of sensors  $x_1$  and  $x_2$  in  $\mathbb{R}^2$  localized by algorithm (1089) as indicated by blue bullets  $\bullet$ . Anchors  $\check{x}_3$ ,  $\check{x}_4$ ,  $\check{x}_5$ , corresponding to Figure 148, are indicated by  $\otimes$ . When targets  $\circ$  and bullets  $\bullet$  coincide under these noiseless conditions, localization is successful. On this run, two visible localization errors are due to rank-3 Gram optimal solutions. These errors can be corrected by choosing a different normal in objective of minimization.

**5.4.2.2.11 Example.** *Tandem trilateration in wireless sensor network.*

Given three known absolute point-positions in  $\mathbb{R}^2$  (three anchors  $\check{x}_3, \check{x}_4, \check{x}_5$ ), two unknown sensors  $x_1, x_2 \in \mathbb{R}^2$  have absolute position determinable from their noiseless distances-square (as indicated in Figure 148) assuming the anchors exhibit no rotational or reflective symmetry in their affine hull (§5.5.2). This example differs from Example 5.4.2.2.8 insofar as trilateration of each sensor is now in terms of one unknown position: the other sensor. We express this localization as a convex optimization problem (a semidefinite program, §4.1) in terms of list  $X \triangleq [x_1 \ x_2 \ \check{x}_3 \ \check{x}_4 \ \check{x}_5] \in \mathbb{R}^{2 \times 5}$  and Gram matrix  $G \in \mathbb{S}^5$  (1042) via relaxation (1082):

$$\begin{aligned} & \underset{G \in \mathbb{S}^5, X \in \mathbb{R}^{2 \times 5}}{\text{minimize}} \quad \text{tr } G \\ & \text{subject to} \quad \begin{aligned} \text{tr}(G\Phi_{i1}) &= \check{d}_{i1}, & i &= 2, 4, 5 \\ \text{tr}(G\Phi_{i2}) &= \check{d}_{i2}, & i &= 3, 5 \\ \text{tr}(Ge_i e_i^T) &= \|\check{x}_i\|^2, & i &= 3, 4, 5 \\ \text{tr}(Ge_i e_j^T + e_j e_i^T)/2 &= \check{x}_i^T \check{x}_j, & 3 \leq i < j &= 4, 5 \\ X(:, 3:5) &= [\check{x}_3 \ \check{x}_4 \ \check{x}_5] \\ \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} &\succeq 0 \end{aligned} \end{aligned} \tag{1089}$$

where

$$\Phi_{ij} = (e_i - e_j)(e_i - e_j)^T \in \mathbb{S}_+^N \tag{1031}$$

This problem realization is fragile because of the unknown distances between sensors and anchors. Yet there is no more information we may include beyond the 11 independent equality constraints on the Gram matrix (nonredundant constraints not antithetical) to reduce the feasible set.<sup>5.17</sup>

Exhibited in Figure 149 are two mistakes in solution  $X^*(:, 1:2)$  due to a rank-3 optimal Gram matrix  $G^*$ . The trace objective is a heuristic minimizing convex envelope of quasiconcave function<sup>5.18</sup>  $\text{rank } G$ . (§2.9.2.9.2, §7.2.2.1) A rank-2 optimal Gram matrix can be found and the errors corrected by choosing a different normal for the linear objective function, now implicitly the Identity matrix  $I$ ; *id est*,

$$\text{tr } G = \langle G, I \rangle \leftarrow \langle G, \delta(u) \rangle \tag{1090}$$

where vector  $u \in \mathbb{R}^5$  is randomly selected. A random search for a good normal  $\delta(u)$  in only a few iterations is quite easy and effective because: the problem is small, an optimal solution is known *a priori* to exist in two dimensions, a good normal direction is not necessarily unique, and (we speculate) because the feasible affine-subset slices the positive semidefinite cone thinly in the Euclidean sense.<sup>5.19</sup>  $\square$

We explore ramifications of noise and incomplete data throughout; their individual effect being to expand the optimal solution set, introducing more solutions and higher-rank solutions. Hence our focus shifts in §4.5 to discovery of a reliable means for diminishing the optimal solution set by introduction of a rank constraint.

Now we illustrate how a problem in distance geometry can be solved without equality constraints representing measured distance; instead, we have only upper and lower bounds on distances measured:

<sup>5.17</sup>By virtue of their dimensioning, the sensors are already constrained to  $\mathbb{R}^2$  the affine hull of the anchors.

<sup>5.18</sup>Projection on that nonconvex subset of all  $N \times N$ -dimensional positive semidefinite matrices, in an affine subset, whose rank does not exceed 2 is a problem considered difficult to solve. [394, §4]

<sup>5.19</sup>The log det rank-heuristic from §7.2.2.4 does not work here because it chooses the wrong normal. Rank reduction (§4.1.2.1) is unsuccessful here because Barvinok's upper bound (§2.9.3.0.1) on rank of  $G^*$  is 4.

#### 5.4.2.2.12 Example. Wireless location in a cellular telephone network.

Utilizing measurements of distance, time of flight, angle of arrival, or signal power in the context of wireless telephony, *multilateration* is the process of localizing (determining absolute position of) a radio signal source • by inferring geometry relative to multiple fixed *base stations* ○ whose locations are known.

We consider localization of a cellular telephone by distance geometry, so we assume distance to any particular base station can be inferred from received signal power. On a large open flat expanse of terrain, signal-power measurement corresponds well with inverse distance. But it is not uncommon for measurement of signal power to suffer 20 decibels in loss caused by factors such as *multipath* interference (signal reflections), mountainous terrain, man-made structures, turning one's head, or rolling the windows up in an automobile. Consequently, contours of equal signal power are no longer circular; their geometry is irregular and would more aptly be approximated by translated ellipsoids of graduated orientation and eccentricity as in Figure 151.

Depicted in Figure 150 is one cell phone  $x_1$  whose signal power is automatically and repeatedly measured by 6 base stations ○ nearby.<sup>5.20</sup> Those signal power measurements are transmitted from that cell phone to base station  $\check{x}_2$  who decides whether to transfer (*hand-off* or *hand-over*) responsibility for that call should the user roam outside its cell.<sup>5.21</sup>

Due to noise, at least one distance measurement more than the minimum number of measurements is required for reliable localization in practice; 3 measurements are minimum in two dimensions, 4 in three.<sup>5.22</sup> Existence of noise precludes measured distance from the input data. We instead assign measured distance to a range estimate specified by individual upper and lower bounds:  $\overline{d_{i1}}$  is the upper bound on distance-square from the cell phone to  $i^{\text{th}}$  base station, while  $\underline{d_{i1}}$  is the lower bound. These bounds become the input data. Each measurement range is presumed different from the others.

Then convex problem (1081) takes the form:

$$\begin{aligned} & \underset{G \in \mathbb{S}^7, X \in \mathbb{R}^{2 \times 7}}{\text{minimize}} \quad \text{tr } G \\ \text{subject to} \quad & \begin{aligned} \underline{d_{i1}} \leq \text{tr}(G \Phi_{i1}) & \leq \overline{d_{i1}}, & i = 2 \dots 7 \\ \text{tr}(G e_i e_i^T) & = \|\check{x}_i\|^2, & i = 2 \dots 7 \\ \text{tr}(G(e_i e_j^T + e_j e_i^T)/2) & = \check{x}_i^T \check{x}_j, & 2 \leq i < j = 3 \dots 7 \\ X(:, 2:7) & = [\check{x}_2 \check{x}_3 \check{x}_4 \check{x}_5 \check{x}_6 \check{x}_7] \\ \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} & \succeq 0 \end{aligned} \end{aligned} \tag{1091}$$

where

$$\Phi_{ij} = (e_i - e_j)(e_i - e_j)^T \in \mathbb{S}_+^N \tag{1031}$$

This semidefinite program realizes the wireless location problem illustrated in Figure 150. Location  $X^*(\cdot, 1)$  is taken as solution, although measurement noise will often cause  $\text{rank } G^*$  to exceed 2. Randomized search for a rank-2 optimal solution is not so easy here as in Example 5.4.2.2.11. We introduce a method in §4.5 for enforcing the stronger rank-constraint (1085). To formulate this same problem in three dimensions, point list  $X$  is simply redimensioned in the semidefinite program. □

<sup>5.20</sup>Cell phone signal power is typically encoded logarithmically with 1-decibel increment and 64-decibel dynamic range.

<sup>5.21</sup>Because distance to base station is quite difficult to infer from signal power measurements in an urban environment, localization of a particular cell phone • by distance geometry would be far easier were the whole cellular system instead conceived so cell phone  $x_1$  also transmits (to base station  $\check{x}_2$ ) its signal power as received by all other cell phones within range.

<sup>5.22</sup>In Example 4.5.1.2.4, we explore how this convex optimization algorithm fares in the face of measurement noise.

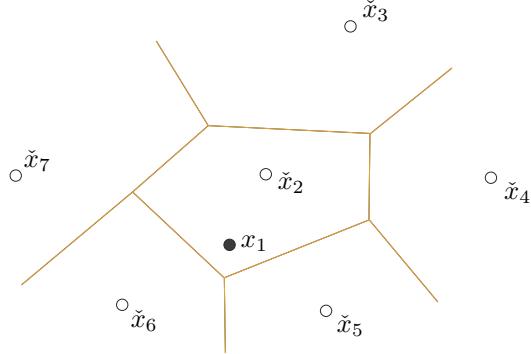


Figure 150: Regions of coverage by base stations  $\circ$  in a cellular telephone network. The term *cellular* arises from packing of regions best covered by neighboring base stations. Illustrated is a pentagonal *cell* best covered by base station  $\check{x}_2$ . Like a Voronoi diagram, cell geometry depends on base-station arrangement. In some US urban environments, it is not unusual to find base stations spaced approximately 1 mile apart. There can be as many as 20 base-station antennae capable of receiving signal from any given cell phone  $\bullet$ ; practically, that number is closer to 6.

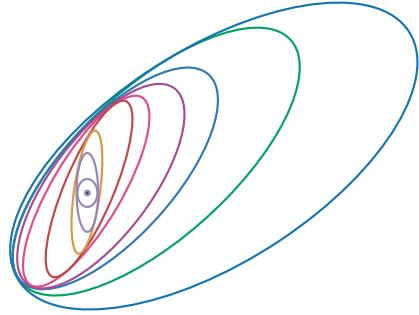


Figure 151: Some fitted contours of equal signal power in  $\mathbb{R}^2$  transmitted from a commercial cellular telephone  $\bullet$  over about 1 mile suburban terrain outside San Francisco in 2005. (Data by courtesy of [Polaris Wireless](#).)

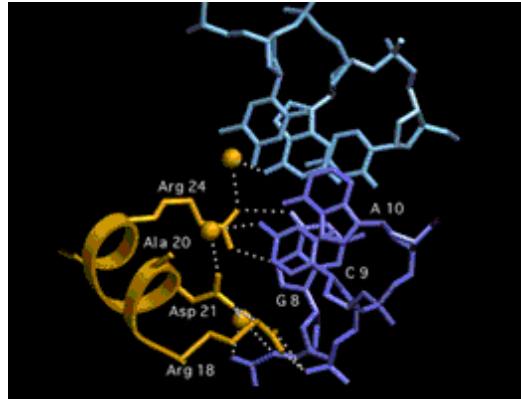


Figure 152: A depiction of molecular conformation. [134]

#### 5.4.2.2.13 Example. (Biswas, Nigam, Ye) Molecular Conformation.

The subatomic measurement technique called *nuclear magnetic resonance spectroscopy* (NMR) is employed to ascertain physical conformation of molecules; *e.g.*, Figure 5, Figure 152. From this technique, distance, angle, and dihedral angle measurements can be obtained. Dihedral angles arise consequent to a phenomenon where atom subsets are physically constrained to Euclidean planes.

*In the rigid covalent geometry approximation, the bond lengths and angles are treated as completely fixed, so that a given spatial structure can be described very compactly indeed by a list of torsion angles alone... These are the dihedral angles between the planes spanned by the two consecutive triples in a chain of four covalently bonded atoms.*

—G. M. Crippen & T. F. Havel, 1988 [96, §1.1]

Crippen & Havel recommend working exclusively with distance data because they consider angle data to be mathematically cumbersome. The present example shows instead how inclusion of dihedral angle data into a problem statement can be made elegant and convex.

As before, ascribe position information to the matrix

$$X = [x_1 \cdots x_N] \in \mathbb{R}^{3 \times N} \quad (77)$$

and introduce a matrix  $\aleph$  holding normals  $\eta$  to planes respecting dihedral angles  $\varphi$ :

$$\aleph \triangleq [\eta_1 \cdots \eta_M] \in \mathbb{R}^{3 \times M} \quad (1092)$$

As in the other examples, we preferentially work with Gram matrices  $G$  because of the bridge they provide between other variables; we define

$$\begin{bmatrix} G_\aleph & Z \\ Z^T & G_X \end{bmatrix} \triangleq \begin{bmatrix} \aleph^T \aleph & \aleph^T X \\ X^T \aleph & X^T X \end{bmatrix} = \begin{bmatrix} \aleph^T \\ X^T \end{bmatrix} [\aleph \ X] \in \mathbb{R}^{N+M \times N+M} \quad (1093)$$

whose rank is 3 by assumption. So our problem's variables are the two Gram matrices  $G_X$  and  $G_\aleph$  and matrix  $Z = \aleph^T X$  of cross products. Then measurements of distance-square  $d$  can be expressed as linear constraints on  $G_X$  as in (1091), dihedral angle  $\varphi$  measurements can be expressed as linear constraints on  $G_\aleph$  by (1061), and normal-vector  $\eta$  conditions can be expressed by vanishing linear constraints on cross-product matrix  $Z$ : Consider

three points  $x$  labelled 1, 2, 3 assumed to lie in the  $\ell^{\text{th}}$  plane whose normal is  $\eta_\ell$ . There might occur, for example, the independent constraints

$$\begin{aligned}\eta_\ell^T(x_1 - x_2) &= 0 \\ \eta_\ell^T(x_2 - x_3) &= 0\end{aligned}\quad (1094)$$

which are expressible in terms of constant matrices  $A_k \in \mathbb{R}^{M \times N}$ ;

$$\begin{aligned}\langle Z, A_{\ell 12} \rangle &= 0 \\ \langle Z, A_{\ell 23} \rangle &= 0\end{aligned}\quad (1095)$$

Although normals  $\eta$  can be constrained exactly to unit length,

$$\delta(G_N) = \mathbf{1} \quad (1096)$$

NMR data is noisy; so measurements are given as upper and lower bounds. Given bounds on dihedral angles respecting  $0 \leq \varphi_j \leq \pi$  and bounds on distances  $d_i$  and given constant matrices  $A_k$  (1095) and symmetric matrices  $\Phi_i$  (1031) and  $B_j$  per (1061), then a molecular conformation problem can be expressed:

$$\begin{array}{ll}\text{find} & G_X \\ \text{subject to} & \begin{array}{lcl}\underline{d}_i & \leq & \text{tr}(G_X \Phi_i) & \leq & \bar{d}_i & \forall i \in \mathcal{I}_1 \\ \cos \varphi_j & \leq & \text{tr}(G_N B_j) & \leq & \cos \varphi_j & \forall j \in \mathcal{I}_2 \\ \langle Z, A_k \rangle & = & 0 & & & \forall k \in \mathcal{I}_3 \\ G_X \mathbf{1} & = & \mathbf{0} & & & \\ \delta(G_N) & = & \mathbf{1} & & & \\ \begin{bmatrix} G_N & Z \\ Z^T & G_X \end{bmatrix} & \succeq & 0 & & & \\ \text{rank} \begin{bmatrix} G_N & Z \\ Z^T & G_X \end{bmatrix} & = & 3 & & & \end{array}\end{array}\quad (1097)$$

where  $G_X \mathbf{1} = \mathbf{0}$  provides a geometrically centered list  $X$  (§5.4.2.2). Ignoring the rank constraint would tend to force cross-product matrix  $Z$  to zero. What binds these three variables is the rank constraint; we show how to satisfy it in §4.5.  $\square$

### 5.4.3 Inner-product form EDM definition

We might, for example, want to realize a constellation given only interstellar distance (or, equivalently, parsecs from our Sun and relative angular measurement; the Sun as vertex to two distant stars); called stellar cartography...  
–p.19

Equivalent to (1029) is [442, §1-7] [368, §3.2]

$$\begin{aligned}d_{ij} &= d_{ik} + d_{kj} - 2\sqrt{d_{ik}d_{kj}} \cos \theta_{ikj} \\ &= [\sqrt{d_{ik}} \quad \sqrt{d_{kj}}] \begin{bmatrix} 1 & -e^{i\theta_{ikj}} \\ -e^{-i\theta_{ikj}} & 1 \end{bmatrix} \begin{bmatrix} \sqrt{d_{ik}} \\ \sqrt{d_{kj}} \end{bmatrix}\end{aligned}\quad (1098)$$

called *law of cosines* where  $i \triangleq \sqrt{-1}$ ,  $i, j, k$  are positive integers, and  $\theta_{ikj}$  is the angle at vertex  $x_k$  formed by vectors  $x_i - x_k$  and  $x_j - x_k$ ; *id est*, the angle relative to  $x_k$

$$\cos \theta_{ikj} = \frac{\frac{1}{2}(d_{ik} + d_{kj} - d_{ij})}{\sqrt{d_{ik}d_{kj}}} = \frac{(x_i - x_k)^T(x_j - x_k)}{\|x_i - x_k\| \|x_j - x_k\|} \quad (1099)$$

where the numerator forms an inner product of vectors. Distance-square  $d_{ij} \left( \begin{bmatrix} \sqrt{d_{ik}} \\ \sqrt{d_{kj}} \end{bmatrix} \right)$  is a convex quadratic function<sup>5.23</sup> on  $\mathbb{R}_+^2$  whereas  $d_{ij}(\theta_{ikj})$  is quasiconvex (§3.14) minimized over domain  $\{-\pi \leq \theta_{ikj} \leq \pi\}$  by  $\theta_{ikj}^* = 0$ , we get the *Pythagorean theorem* when  $\theta_{ikj} = \pm\pi/2$ , and  $d_{ij}(\theta_{ikj})$  is maximized when  $\theta_{ikj}^* = \pm\pi$ ;

$$\begin{aligned} d_{ij} &= (\sqrt{d_{ik}} + \sqrt{d_{kj}})^2, \quad \theta_{ikj} = \pm\pi \\ d_{ij} &= d_{ik} + d_{kj}, \quad \theta_{ikj} = \pm\frac{\pi}{2} \\ d_{ij} &= (\sqrt{d_{ik}} - \sqrt{d_{kj}})^2, \quad \theta_{ikj} = 0 \end{aligned} \quad (1100)$$

so

$$|\sqrt{d_{ik}} - \sqrt{d_{kj}}| \leq \sqrt{d_{ij}} \leq \sqrt{d_{ik}} + \sqrt{d_{kj}} \quad (1101)$$

Hence the triangle inequality, Euclidean metric property 4, holds for any EDM  $D$ .

We may construct an inner-product form of the EDM definition for matrices by evaluating (1098) for  $k=1$ : By defining

$$\Theta^T \Theta \triangleq \begin{bmatrix} d_{12} & \sqrt{d_{12}d_{13}} \cos \theta_{213} & \sqrt{d_{12}d_{14}} \cos \theta_{214} & \cdots & \sqrt{d_{12}d_{1N}} \cos \theta_{21N} \\ \sqrt{d_{12}d_{13}} \cos \theta_{213} & d_{13} & \sqrt{d_{13}d_{14}} \cos \theta_{314} & \cdots & \sqrt{d_{13}d_{1N}} \cos \theta_{31N} \\ \sqrt{d_{12}d_{14}} \cos \theta_{214} & \sqrt{d_{13}d_{14}} \cos \theta_{314} & d_{14} & \ddots & \sqrt{d_{14}d_{1N}} \cos \theta_{41N} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \sqrt{d_{12}d_{1N}} \cos \theta_{21N} & \sqrt{d_{13}d_{1N}} \cos \theta_{31N} & \sqrt{d_{14}d_{1N}} \cos \theta_{41N} & \cdots & d_{1N} \end{bmatrix} \in \mathbb{S}^{N-1} \quad (1102)$$

then any EDM may be expressed

$$\begin{aligned} \mathbf{D}(\Theta) &\triangleq \begin{bmatrix} 0 \\ \delta(\Theta^T \Theta) \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & \delta(\Theta^T \Theta)^T \end{bmatrix} - 2 \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \Theta^T \Theta \end{bmatrix} \in \mathbb{EDM}^N \\ &= \begin{bmatrix} 0 & \delta(\Theta^T \Theta)^T \\ \delta(\Theta^T \Theta) & \delta(\Theta^T \Theta) \mathbf{1}^T + \mathbf{1} \delta(\Theta^T \Theta)^T - 2 \Theta^T \Theta \end{bmatrix} \end{aligned} \quad (1103)$$

$$\mathbb{EDM}^N = \left\{ \mathbf{D}(\Theta) \mid \Theta \in \mathbb{R}^{N-1 \times N-1} \right\} \quad (1104)$$

for which all Euclidean metric properties hold. Entries of  $\Theta^T \Theta$  result from vector inner-products as in (1099); *id est*,

$$\Theta = [x_2 - x_1 \quad x_3 - x_1 \quad \cdots \quad x_N - x_1] = X\sqrt{2}V_N \in \mathbb{R}^{n \times N-1} \quad (1105)$$

Inner product  $\Theta^T \Theta$  is obviously related to a Gram matrix (1042),

$$G = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \Theta^T \Theta \end{bmatrix}, \quad x_1 = \mathbf{0} \quad (1106)$$

For  $D = \mathbf{D}(\Theta)$  and no condition on the list  $X$  (*confer*(1050)(1054))

$$\Theta^T \Theta = -V_N^T D V_N \in \mathbb{R}^{N-1 \times N-1} \quad (1107)$$

---

<sup>5.23</sup>  $\begin{bmatrix} 1 & -e^{i\theta_{ikj}} \\ -e^{-i\theta_{ikj}} & 1 \end{bmatrix} \succeq 0$ , having eigenvalues  $\{0, 2\}$ .

Minimum is attained for  $\begin{bmatrix} \sqrt{d_{ik}} \\ \sqrt{d_{kj}} \end{bmatrix} = \begin{cases} \mu \mathbf{1}, & \mu \geq 0, \theta_{ikj} = 0 \\ \mathbf{0}, & -\pi \leq \theta_{ikj} \leq \pi, \theta_{ikj} \neq 0 \end{cases}$  (§D.2.1, [65, exmp.4.5]).

### 5.4.3.1 Relative-angle form

The inner-product form EDM definition is not a unique definition of Euclidean distance matrix; there are approximately five flavors distinguished by their argument to operator  $\mathbf{D}$ . Here is another one:

Like  $\mathbf{D}(X)$  (1033),  $\mathbf{D}(\Theta)$  will make an EDM given any  $\Theta \in \mathbb{R}^{n \times N-1}$ , it is neither a convex function of  $\Theta$  (§5.4.3.2), and it is homogeneous in the sense (1036). Scrutinizing  $\Theta^T \Theta$  (1102) we find that because of the arbitrary choice  $k=1$ , distances therein are all with respect to point  $x_1$ . Similarly, relative angles in  $\Theta^T \Theta$  are between all vector pairs having vertex  $x_1$ . Yet picking arbitrary  $\theta_{i1j}$  to fill  $\Theta^T \Theta$  will not necessarily make an EDM; inner product (1102) must be positive semidefinite.

$$\Theta^T \Theta = \sqrt{\delta(d)} \Omega \sqrt{\delta(d)} \triangleq$$

$$\begin{bmatrix} \sqrt{d_{12}} & & \mathbf{0} \\ & \sqrt{d_{13}} & \\ & \ddots & \\ \mathbf{0} & & \sqrt{d_{1N}} \end{bmatrix} \begin{bmatrix} 1 & \cos \theta_{213} & \cdots & \cos \theta_{21N} \\ \cos \theta_{213} & 1 & \ddots & \cos \theta_{31N} \\ \vdots & \ddots & \ddots & \vdots \\ \cos \theta_{21N} & \cos \theta_{31N} & \cdots & 1 \end{bmatrix} \begin{bmatrix} \sqrt{d_{12}} & & \mathbf{0} \\ & \sqrt{d_{13}} & \\ & \ddots & \\ \mathbf{0} & & \sqrt{d_{1N}} \end{bmatrix} \quad (1108)$$

Expression  $\mathbf{D}(\Theta)$  defines an EDM for any positive semidefinite *relative-angle matrix*

$$\Omega = [\cos \theta_{i1j}, i, j = 2 \dots N] \in \mathbb{S}^{N-1} \quad (1109)$$

and any nonnegative distance vector

$$d = [d_{1j}, j = 2 \dots N] = \delta(\Theta^T \Theta) \in \mathbb{R}^{N-1} \quad (1110)$$

because (§A.3.1.0.5)

$$\Omega \succeq 0 \Rightarrow \Theta^T \Theta \succeq 0 \quad (1111)$$

Decomposition (1108) and the *relative-angle matrix inequality*  $\Omega \succeq 0$  lead to a different expression of an inner-product form EDM definition (1103)

$$\begin{aligned} \mathbf{D}(\Omega, d) &\triangleq \begin{bmatrix} 0 \\ d \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & d^T \end{bmatrix} - 2\sqrt{\delta\left(\begin{bmatrix} 0 \\ d \end{bmatrix}\right)} \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \Omega \end{bmatrix} \sqrt{\delta\left(\begin{bmatrix} 0 \\ d \end{bmatrix}\right)} \\ &= \begin{bmatrix} 0 & d^T \\ d & d\mathbf{1}^T + \mathbf{1}d^T - 2\sqrt{\delta(d)} \Omega \sqrt{\delta(d)} \end{bmatrix} \in \mathbb{EDM}^N \end{aligned} \quad (1112)$$

and another expression of the EDM cone:

$$\mathbb{EDM}^N = \left\{ \mathbf{D}(\Omega, d) \mid \Omega \succeq 0, \sqrt{\delta(d)} \succeq 0 \right\} \quad (1113)$$

In the particular circumstance  $x_1 = \mathbf{0}$ , we can relate interpoint angle matrix  $\Psi$  from the Gram decomposition in (1042) to relative-angle matrix  $\Omega$  in (1108). Thus,

$$\Psi \equiv \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \Omega \end{bmatrix}, \quad x_1 = \mathbf{0} \quad (1114)$$

### 5.4.3.2 Inner-product form $-V_N^T \mathbf{D}(\Theta) V_N$ convexity

On page 368 we saw that each EDM entry  $d_{ij}$  is a convex quadratic function of  $\begin{bmatrix} \sqrt{d_{ik}} \\ \sqrt{d_{kj}} \end{bmatrix}$  and a quasiconvex function of  $\theta_{ikj}$ . Here the situation for inner-product form EDM

operator  $\mathbf{D}(\Theta)$  (1103) is identical to that in §5.4.1 for list-form  $\mathbf{D}(X)$ ;  $-\mathbf{D}(\Theta)$  is not a quasiconvex function of  $\Theta$  by the same reasoning, and from (1107)

$$-V_N^T \mathbf{D}(\Theta) V_N = \Theta^T \Theta \quad (1115)$$

is a convex quadratic function of  $\Theta$  on domain  $\mathbb{R}^{n \times N-1}$  achieving its minimum at  $\Theta = \mathbf{0}$ .

### 5.4.3.3 Inner-product form, discussion

We deduce that knowledge of interpoint distance is equivalent to knowledge of distance and angle from the perspective of one point,  $x_1$  in our chosen case. The total amount of information  $N(N-1)/2$  in  $\Theta^T \Theta$  is unchanged<sup>5.24</sup> with respect to EDM  $D$ .

## 5.5 Invariance

When  $D$  is an EDM, there exist an infinite number of corresponding  $N$ -point lists  $X$  (77) in Euclidean space. All those lists are related by isometric transformation: rotation, reflection, and translation (*offset* or *shift*).

### 5.5.1 Translation

Any translation common among all the points  $x_\ell$  in a list will be cancelled in the formation of each  $d_{ij}$ . Proof follows directly from (1029). Knowing that translation  $\alpha$  in advance, we may remove it from the list constituting the columns of  $X$  by subtracting  $\alpha \mathbf{1}^T$ . Then it stands to reason by list-form definition (1033) of an EDM, for any translation  $\alpha \in \mathbb{R}^n$

$$\mathbf{D}(X - \alpha \mathbf{1}^T) = \mathbf{D}(X) \quad (1116)$$

In words, interpoint distances are unaffected by offset; EDM  $D$  is *translation invariant*. When  $\alpha = x_1$  in particular,

$$[x_2 - x_1 \ x_3 - x_1 \ \cdots \ x_N - x_1] = X\sqrt{2}V_N \in \mathbb{R}^{n \times N-1} \quad (1105)$$

and so

$$\mathbf{D}(X - x_1 \mathbf{1}^T) = \mathbf{D}(X - X e_1 \mathbf{1}^T) = \mathbf{D}\left(X \begin{bmatrix} \mathbf{0} & \sqrt{2}V_N \end{bmatrix}\right) = \mathbf{D}(X) \quad (1117)$$

#### 5.5.1.0.1 Example. Translating geometric center to origin.

We might choose to shift the geometric center  $\alpha_c$  of an  $N$ -point list  $\{x_\ell\}$  (arranged columnar in  $X$ ) to the origin; [393] [187]

$$\alpha = \alpha_c \triangleq X b_c \triangleq X \mathbf{1} \frac{1}{N} \in \mathcal{P} \subseteq \mathcal{A} \quad (1118)$$

where  $\mathcal{A}$  represents the list's affine hull. If we were to associate a point-mass  $m_\ell$  with each of the points  $x_\ell$  in the list, then their *center of mass* (or *gravity*) would be  $(\sum x_\ell m_\ell) / \sum m_\ell$ . The geometric center is the same as the center of mass under the assumption of uniform mass density across points. [387] The geometric center always lies in

---

<sup>5.24</sup>The reason for amount  $O(N^2)$  information is because of the relative measurements. Use of a fixed reference in measurement of angles and distances would reduce required information but is antithetical. In the particular case  $n=2$ , for example, ordering all points  $x_\ell$  (in a length- $N$  list) by increasing angle of vector  $x_\ell - x_1$  with respect to  $x_2 - x_1$ ,  $\theta_{i1j}$  becomes equivalent to  $\sum_{k=i}^{j-1} \theta_{k,1,k+1} \leq 2\pi$  and the amount of information is reduced to  $2N-3$ ; rather,  $O(N)$ .

the convex hull  $\mathcal{P}$  of the list; *id est*,  $\alpha_c \in \mathcal{P}$  because  $b_c^T \mathbf{1} = 1$  and  $b_c \succeq 0$ .<sup>5.25</sup> Subtracting the geometric center from every list member,

$$X - \alpha_c \mathbf{1}^T = X - \frac{1}{N} X \mathbf{1} \mathbf{1}^T = X(I - \frac{1}{N} \mathbf{1} \mathbf{1}^T) = XV \in \mathbb{R}^{n \times N} \quad (1119)$$

where  $V$  is the geometric centering matrix (1055). So we have (*confer* (1033))

$$\mathbf{D}(X) = \mathbf{D}(XV) = \delta(V^T X^T XV) \mathbf{1}^T + \mathbf{1} \delta(V^T X^T XV)^T - 2V^T X^T XV \in \mathbb{EDM}^N \quad (1120)$$

□

### 5.5.1.1 Gram-form invariance

Following from (1120) and the linear Gram-form EDM operator (1045):

$$\mathbf{D}(G) = \mathbf{D}(VGV) = \delta(VGV) \mathbf{1}^T + \mathbf{1} \delta(VGV)^T - 2VGV \in \mathbb{EDM}^N \quad (1121)$$

The Gram-form consequently exhibits invariance to translation by a *doublet*  $u\mathbf{1}^T + \mathbf{1}u^T$  ([§B.2](#))

$$\mathbf{D}(G) = \mathbf{D}(G - (u\mathbf{1}^T + \mathbf{1}u^T)) \quad (1122)$$

because, for any  $u \in \mathbb{R}^N$ ,  $\mathbf{D}(u\mathbf{1}^T + \mathbf{1}u^T) = \mathbf{0}$ . The collection of all such doublets forms the nullspace (1138) to the operator; the *translation-invariant subspace*  $\mathbb{S}_c^{N\perp}$  (2198) of the symmetric matrices  $\mathbb{S}^N$ . This means matrix  $G$  is not unique and can belong to an expanse more broad than a positive semidefinite cone; *id est*,  $G \in \mathbb{S}_+^N - \mathbb{S}_c^{N\perp}$ . So explains Gram matrix sufficiency in EDM definition (1045).<sup>5.26</sup>

### 5.5.2 Rotation/Reflection

Rotation of the list  $X \in \mathbb{R}^{n \times N}$  about some arbitrary point  $\alpha \in \mathbb{R}^n$ , or reflection through some affine subset containing  $\alpha$ , can be accomplished via  $Q(X - \alpha \mathbf{1}^T)$  where  $Q$  is an orthogonal matrix ([§B.5](#)).

We rightfully expect

$$\mathbf{D}(Q(X - \alpha \mathbf{1}^T)) = \mathbf{D}(QX - \beta \mathbf{1}^T) = \mathbf{D}(QX) = \mathbf{D}(X) \quad (1123)$$

Because list-form  $\mathbf{D}(X)$  is translation invariant, we may safely ignore offset and consider only the impact of matrices that premultiply  $X$ . Interpoint distances are unaffected by rotation or reflection; we say, EDM  $D$  is *rotation/reflection invariant*. Proof follows from the fact,  $Q^T = Q^{-1} \Rightarrow X^T Q^T Q X = X^T X$ . So (1123) follows directly from (1033).

The class of premultiplying matrices for which interpoint distances are unaffected is a little more broad than orthogonal matrices. Looking at EDM definition (1033), it appears that any matrix  $Q_p$  such that

$$X^T Q_p^T Q_p X = X^T X \quad (1124)$$

will have the property

$$\mathbf{D}(Q_p X) = \mathbf{D}(X) \quad (1125)$$

An example is thin  $Q_p \in \mathbb{R}^{m \times n}$  ( $m > n$ ) having orthonormal columns; an orthonormal matrix.

---

<sup>5.25</sup> Any  $b$  from  $\alpha = Xb$  chosen such that  $b^T \mathbf{1} = 1$ , more generally, makes an auxiliary  $V$ -matrix. ([§B.4.5](#))

<sup>5.26</sup> A constraint  $G\mathbf{1} = \mathbf{0}$  would prevent excursion into the translation-invariant subspace (numerical unboundedness).

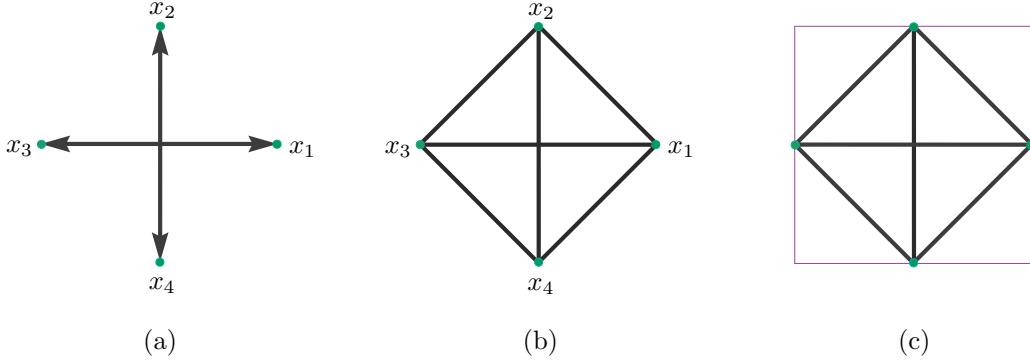


Figure 153: (a) Four points in quadrature in two dimensions about their geometric center. (b) Complete EDM graph of diamond-shaped vertices. (c) Quadrature rotation of Euclidean body in  $\mathbb{R}^2$  first requires shroud: the smallest Cartesian square containing it.

#### 5.5.2.0.1 Example. Reflection prevention and quadrature rotation.

Consider the EDM graph in Figure 153b representing known distance between vertices (Figure 153a) of a tilted-square diamond in  $\mathbb{R}^2$ . Suppose some geometrical optimization problem were posed where isometric transformation is allowed excepting reflection, and where rotation must be quantized so that only *quadrature* rotations are allowed; only multiples of  $\pi/2$ .

In two dimensions, a counterclockwise rotation of any vector about the origin by angle  $\theta$  is prescribed by the orthogonal matrix

$$Q = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad (1126)$$

whereas reflection of any point through a hyperplane containing the origin

$$\partial\mathcal{H} = \left\{ x \in \mathbb{R}^2 \mid \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}^T x = 0 \right\} \quad (1127)$$

is accomplished via multiplication with symmetric orthogonal matrix (§B.5.3)

$$R = \begin{bmatrix} \sin(\theta)^2 - \cos(\theta)^2 & -2\sin(\theta)\cos(\theta) \\ -2\sin(\theta)\cos(\theta) & \cos(\theta)^2 - \sin(\theta)^2 \end{bmatrix} \quad (1128)$$

Rotation matrix  $Q$  is characterized by identical diagonal entries and by antidiagonal entries equal but opposite in sign, whereas reflection matrix  $R$  is characterized in the reverse sense.

Assign the diamond vertices  $\{x_\ell \in \mathbb{R}^2, \ell=1 \dots 4\}$  to columns of a matrix

$$X = [x_1 \ x_2 \ x_3 \ x_4] \in \mathbb{R}^{2 \times 4} \quad (77)$$

Our scheme to prevent reflection enforces a rotation matrix characteristic upon the coordinates of adjacent points themselves: First shift the geometric center of  $X$  to the origin; for geometric centering matrix  $V \in \mathbb{S}^4$  (§5.5.1.0.1), define

$$Y \triangleq X V \in \mathbb{R}^{2 \times 4} \quad (1129)$$

To maintain relative quadrature between points (Figure 153a) and to prevent reflection, it is sufficient that all interpoint distances be specified and that adjacencies  $Y(:, 1:2)$ ,

$Y(:, 2:3)$ , and  $Y(:, 3:4)$  be proportional to  $2 \times 2$  rotation matrices; any clockwise rotation would ascribe a reflection matrix characteristic. Counterclockwise rotation is thereby enforced by constraining equality among diagonal and antidiagonal entries as prescribed by (1126);

$$Y(:, 1:3) = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} Y(:, 2:4) \quad (1130)$$

Quadrature quantization of rotation can be regarded as a constraint on tilt of the smallest Cartesian square containing the diamond as in Figure 153c. Our scheme to quantize rotation requires that all square vertices be described by vectors whose entries are nonnegative when the square is translated anywhere interior to the nonnegative orthant. We capture the four square vertices as columns of a product  $YC$  where

$$C = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad (1131)$$

Then, assuming a unit-square shroud, the affine constraint

$$YC + \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix} \mathbf{1}^T \geq \mathbf{0} \quad (1132)$$

quantizes rotation, as desired.  $\square$

#### 5.5.2.1 Inner-product form invariance

Likewise,  $\mathbf{D}(\Theta)$  (1103) is rotation/reflection invariant;

$$\mathbf{D}(Q_p\Theta) = \mathbf{D}(Q\Theta) = \mathbf{D}(\Theta) \quad (1133)$$

so (1124) and (1125) similarly apply.

#### 5.5.3 Invariance conclusion

In the making of an EDM, absolute rotation, reflection, and translation information is lost. Given an EDM, reconstruction of point position (§5.12, the list  $X$ ) can be guaranteed correct only in affine dimension  $r$  and relative position. Given a noiseless complete EDM, this isometric reconstruction is unique insofar as every realization of a corresponding list  $X$  is *congruent*:

## 5.6 Injectivity of $\mathbf{D}$ & unique reconstruction

Injectivity implies uniqueness of isometric reconstruction (§5.4.2.2.10); hence, we endeavor to demonstrate it.

EDM operators list-form  $\mathbf{D}(X)$  (1033), Gram-form  $\mathbf{D}(G)$  (1045), and inner-product form  $\mathbf{D}(\Theta)$  (1103) are many-to-one surjections (§5.5) onto the same range; the EDM cone (§6): (*confer* (1046) (1140))

$$\begin{aligned} \text{EDM}^N &= \left\{ \mathbf{D}(X) : \mathbb{R}^{N-1 \times N} \rightarrow \mathbb{S}_h^N \mid X \in \mathbb{R}^{N-1 \times N} \right\} \\ &= \left\{ \mathbf{D}(G) : \mathbb{S}^N \rightarrow \mathbb{S}_h^N \mid G \in \mathbb{S}_+^N - \mathbb{S}_c^{N \perp} \right\} \\ &= \left\{ \mathbf{D}(\Theta) : \mathbb{R}^{N-1 \times N-1} \rightarrow \mathbb{S}_h^N \mid \Theta \in \mathbb{R}^{N-1 \times N-1} \right\} \end{aligned} \quad (1134)$$

where (§5.5.1.1)

$$\mathbb{S}_c^{N \perp} = \{u\mathbf{1}^T + \mathbf{1}u^T \mid u \in \mathbb{R}^N\} \subseteq \mathbb{S}^N \quad (2198)$$

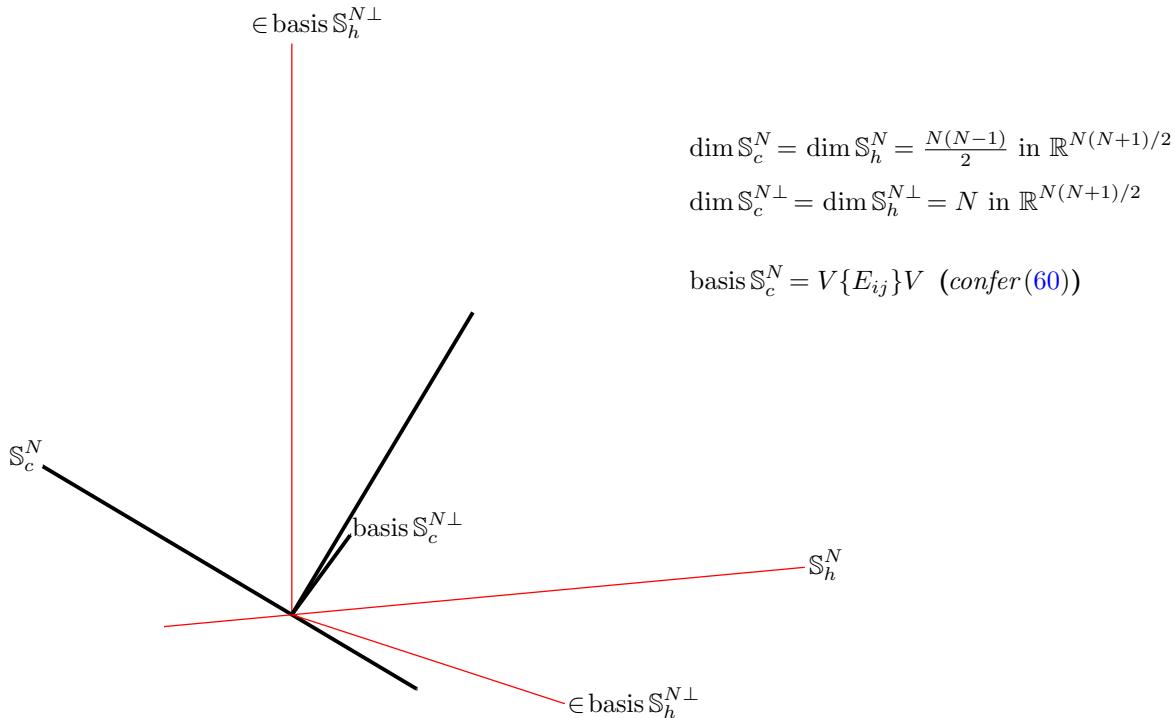


Figure 154: Orthogonal complements in  $\mathbb{S}^N$  abstractly oriented in isometrically isomorphic  $\mathbb{R}^{N(N+1)/2}$ . Case  $N=2$  accurately illustrated in  $\mathbb{R}^3$ . Orthogonal projection of basis for  $\mathbb{S}_h^{N\perp}$  on  $\mathbb{S}_c^{N\perp}$  yields another basis for  $\mathbb{S}_c^{N\perp}$ . (Basis vectors for  $\mathbb{S}_c^{N\perp}$  are illustrated lying in a plane orthogonal to  $\mathbb{S}_c^N$  in this dimension. Basis vectors for each  $\perp$  space outnumber those for its respective orthogonal complement; such is not the case in higher dimension.)

### 5.6.1 Gram-form bijectivity

Because linear Gram-form EDM operator

$$\mathbf{D}(G) = \delta(G)\mathbf{1}^T + \mathbf{1}\delta(G)^T - 2G \quad (1045)$$

has no nullspace [93, §A.1] on the geometric center subspace<sup>5.27</sup> (§E.7.2.0.2)

$$\mathbb{S}_c^N \triangleq \{G \in \mathbb{S}^N \mid G\mathbf{1} = \mathbf{0}\} \quad (2196)$$

$$= \{G \in \mathbb{S}^N \mid \mathcal{N}(G) \supseteq \mathbf{1}\} = \{G \in \mathbb{S}^N \mid \mathcal{R}(G) \subseteq \mathcal{N}(\mathbf{1}^T)\} \quad (1135)$$

$$= \{VYV \mid Y \in \mathbb{S}^N\} \subset \mathbb{S}^N \quad (2197)$$

$$\equiv \{V_{\mathcal{N}}AV_{\mathcal{N}}^T \mid A \in \mathbb{S}^{N-1}\}$$

then  $\mathbf{D}(G)$  on that subspace is injective.

To prove injectivity of  $\mathbf{D}(G)$  on  $\mathbb{S}_c^N$ : Any matrix  $Y \in \mathbb{S}^N$  can be decomposed into orthogonal components in  $\mathbb{S}^N$ :

$$Y = VYV + (Y - VYV) \quad (1136)$$

where  $VYV \in \mathbb{S}_c^N$  and  $Y - VYV \in \mathbb{S}_c^{N\perp}$  (2198). Because of translation invariance (§5.5.1.1) and linearity,  $\mathbf{D}(Y - VYV) = \mathbf{0}$  hence  $\mathcal{N}(\mathbf{D}) \supseteq \mathbb{S}_c^{N\perp}$ . It remains only to show

$$\mathbf{D}(VYV) = \mathbf{0} \Leftrightarrow VYV = \mathbf{0} \quad (1137)$$

( $\Leftrightarrow Y = u\mathbf{1}^T + \mathbf{1}u^T$  for some  $u \in \mathbb{R}^N$ ).  $\mathbf{D}(VYV)$  will vanish whenever  $2VYV = \delta(VYV)\mathbf{1}^T + \mathbf{1}\delta(VYV)^T$ . But this implies  $\mathcal{R}(\mathbf{1})$  (§B.2) were a subset of  $\mathcal{R}(VYV)$ , which is contradictory. Thus we have

$$\mathcal{N}(\mathbf{D}) = \{Y \mid \mathbf{D}(Y) = \mathbf{0}\} = \{Y \mid VYV = \mathbf{0}\} = \mathbb{S}_c^{N\perp} \quad (1138)$$

♦

Since  $G\mathbf{1} = \mathbf{0} \Leftrightarrow X\mathbf{1} = \mathbf{0}$  (1053) simply means list  $X$  is geometrically centered at the origin, and because the Gram-form EDM operator  $\mathbf{D}$  is translation invariant and  $\mathcal{N}(\mathbf{D})$  is the translation-invariant subspace  $\mathbb{S}_c^{N\perp}$ , then EDM definition  $\mathbf{D}(G)$  (1134) on<sup>5.28</sup> (*confer* §6.5.1, §6.6.1, §A.7.4.0.1)

$$\mathbb{S}_c^N \cap \mathbb{S}_+^N = \{VYV \succeq 0 \mid Y \in \mathbb{S}^N\} \equiv \{V_{\mathcal{N}}AV_{\mathcal{N}}^T \mid A \in \mathbb{S}_+^{N-1}\} \subset \mathbb{S}^N \quad (1139)$$

must be surjective onto  $\mathbb{EDM}^N$ ; (*confer* (1046))

$$\mathbb{EDM}^N = \left\{ \mathbf{D}(G) \mid G \in \mathbb{S}_c^N \cap \mathbb{S}_+^N \right\} \quad (1140)$$

<sup>5.27</sup>Equivalence  $\equiv$  in (1135) follows from the fact: Given  $B = VYV = V_{\mathcal{N}}AV_{\mathcal{N}}^T \in \mathbb{S}_c^N$  with only matrix  $A \in \mathbb{S}^{N-1}$  unknown, then  $V_{\mathcal{N}}^\dagger BV_{\mathcal{N}}^{\dagger T} = A$  or  $V_{\mathcal{N}}^\dagger YV_{\mathcal{N}}^{\dagger T} = A$ .

<sup>5.28</sup>Equivalence  $\equiv$  in (1139) follows from the fact: Given  $B = VYV = V_{\mathcal{N}}AV_{\mathcal{N}}^T \in \mathbb{S}_+^N$  with only matrix  $A$  unknown, then  $V_{\mathcal{N}}^\dagger BV_{\mathcal{N}}^{\dagger T} = A$  and  $A \in \mathbb{S}_+^{N-1}$  must be positive semidefinite by positive semidefiniteness of  $B$  and Corollary A.3.1.0.5.

### 5.6.1.1 Gram-form operator $\mathbf{D}$ inversion

Define the linear *geometric centering operator*  $\mathbf{V}$ : (confer (1054))

$$\mathbf{V}(D) : \mathbb{S}^N \rightarrow \mathbb{S}^N \triangleq -VDV^{\frac{1}{2}} \quad (1141)$$

[97, §4.3]<sup>5.29</sup> This orthogonal projector  $\mathbf{V}$  has no nullspace on

$$\mathbb{S}_h^N = \text{aff } \mathbb{EDM}^N \quad (1395)$$

because the projection of  $-D/2$  on  $\mathbb{S}_c^N$  (2196) can be  $\mathbf{0}$  if and only if  $D \in \mathbb{S}_c^{N\perp}$ ; but  $\mathbb{S}_c^{N\perp} \cap \mathbb{S}_h^N = \mathbf{0}$  (Figure 154). Projector  $\mathbf{V}$  on  $\mathbb{S}_h^N$  is therefore injective hence uniquely invertible. Further,  $-V\mathbb{S}_h^NV/2$  is equivalent to the geometric center subspace  $\mathbb{S}_c^N$  in the ambient space of symmetric matrices; a surjection,

$$\mathbb{S}_c^N = \mathbf{V}(\mathbb{S}^N) = \mathbf{V}\left(\mathbb{S}_h^N \oplus \mathbb{S}_h^{N\perp}\right) = \mathbf{V}\left(\mathbb{S}_h^N\right) \quad (1142)$$

because (73)

$$\mathbf{V}\left(\mathbb{S}_h^N\right) \supseteq \mathbf{V}\left(\mathbb{S}_h^{N\perp}\right) = \mathbf{V}\left(\delta^2(\mathbb{S}^N)\right) \quad (1143)$$

Because  $\mathbf{D}(G)$  on  $\mathbb{S}_c^N$  is injective, and  $\text{aff } \mathbf{D}\left(\mathbf{V}(\mathbb{EDM}^N)\right) = \mathbf{D}\left(\mathbf{V}(\text{aff } \mathbb{EDM}^N)\right)$  by property (129) of the affine hull, we find for  $D \in \mathbb{S}_h^N$

$$\mathbf{D}(-VDV^{\frac{1}{2}}) = \delta(-VDV^{\frac{1}{2}})\mathbf{1}^T + \mathbf{1}\delta(-VDV^{\frac{1}{2}})^T - 2(-VDV^{\frac{1}{2}}) \quad (1144)$$

*id est*,

$$D = \mathbf{D}\left(\mathbf{V}(D)\right) \quad (1145)$$

$$-VDV = \mathbf{V}\left(\mathbf{D}(-VDV)\right) \quad (1146)$$

or

$$\mathbb{S}_h^N = \mathbf{D}\left(\mathbf{V}(\mathbb{S}_h^N)\right) \quad (1147)$$

$$-V\mathbb{S}_h^NV = \mathbf{V}\left(\mathbf{D}(-V\mathbb{S}_h^NV)\right) \quad (1148)$$

These operators  $\mathbf{V}$  and  $\mathbf{D}$  are mutual inverses.

The Gram-form  $\mathbf{D}\left(\mathbb{S}_c^N\right)$  (1045) is equivalent to  $\mathbb{S}_h^N$ ;

$$\mathbf{D}\left(\mathbb{S}_c^N\right) = \mathbf{D}\left(\mathbf{V}(\mathbb{S}_h^N \oplus \mathbb{S}_h^{N\perp})\right) = \mathbb{S}_h^N + \mathbf{D}\left(\mathbf{V}(\mathbb{S}_h^{N\perp})\right) = \mathbb{S}_h^N \quad (1149)$$

because  $\mathbb{S}_h^N \supseteq \mathbf{D}\left(\mathbf{V}(\mathbb{S}_h^{N\perp})\right)$ . In summary, for the Gram-form we have the isomorphisms [98, §2] [97, p.76, p.107] [8, §2.1]<sup>5.30</sup> [7, §2] [9, §18.2.1] [3, §2.1]

$$\mathbb{S}_h^N = \mathbf{D}(\mathbb{S}_c^N) \quad (1150)$$

$$\mathbb{S}_c^N = \mathbf{V}(\mathbb{S}_h^N) \quad (1151)$$

and from bijectivity results in §5.6.1,

$$\mathbb{EDM}^N = \mathbf{D}(\mathbb{S}_c^N \cap \mathbb{S}_+^N) \quad (1152)$$

$$\mathbb{S}_c^N \cap \mathbb{S}_+^N = \mathbf{V}(\mathbb{EDM}^N) \quad (1153)$$

<sup>5.29</sup>Critchley cites Torgerson, 1958 [390, ch.11, §2], for a history and derivation of (1141).

<sup>5.30</sup>In [8, p.6, line 20], delete sentence: *Since  $G$  is also... not a singleton set.*

[8, p.10, line 11]  $x_3 = 2$  (not 1).

### 5.6.2 Inner-product form bijectivity

The Gram-form EDM operator  $\mathbf{D}(G) = \delta(G)\mathbf{1}^T + \mathbf{1}\delta(G)^T - 2G$  (1045) is an injective map, for example, on the domain that is the subspace of symmetric matrices having all zeros in the first row and column

$$\begin{aligned}\mathbb{S}_0^N &= \{G \in \mathbb{S}^N \mid Ge_1 = \mathbf{0}\} \\ &= \left\{ \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} Y \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} \mid Y \in \mathbb{S}^N \right\}\end{aligned}\quad (2200)$$

because it obviously has no nullspace there. Since  $Ge_1 = \mathbf{0} \Leftrightarrow Xe_1 = \mathbf{0}$  (1047) means the first point in the list  $X$  resides at the origin, then  $\mathbf{D}(G)$  on  $\mathbb{S}_0^N \cap \mathbb{S}_+^N$  must be surjective onto  $\mathbb{EDM}^N$ .

Substituting  $\Theta^T\Theta \leftarrow -V_N^T DV_N$  (1115) into inner-product form EDM definition  $\mathbf{D}(\Theta)$  (1103), it may be further decomposed:

$$\mathbf{D}(D) = \begin{bmatrix} 0 \\ \delta(-V_N^T DV_N) \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & \delta(-V_N^T DV_N)^T \end{bmatrix} - 2 \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & -V_N^T DV_N \end{bmatrix} \quad (1154)$$

This linear operator  $\mathbf{D}$  is another flavor of inner-product form and an injective map of the EDM cone onto itself. Yet when its domain is instead the entire symmetric hollow subspace  $\mathbb{S}_h^N = \text{aff } \mathbb{EDM}^N$ ,  $\mathbf{D}(D)$  becomes an injective map onto that same subspace. Proof follows directly from the fact: linear  $\mathbf{D}$  has no nullspace [93, §A.1] on  $\mathbb{S}_h^N = \text{aff } \mathbf{D}(\mathbb{EDM}^N) = \mathbf{D}(\text{aff } \mathbb{EDM}^N)$  (129).

#### 5.6.2.1 Inversion of $\mathbf{D}(-V_N^T DV_N)$

Injectivity of  $\mathbf{D}(D)$  suggests inversion of (*confer*(1050))

$$\mathbf{V}_N(D) : \mathbb{S}^N \rightarrow \mathbb{S}^{N-1} \triangleq -V_N^T DV_N \quad (1155)$$

a linear surjective<sup>5.31</sup> mapping onto  $\mathbb{S}^{N-1}$  having nullspace<sup>5.32</sup>  $\mathbb{S}_c^{N\perp}$ ;

$$\mathbf{V}_N(\mathbb{S}_h^N) = \mathbb{S}^{N-1} \quad (1156)$$

injective on domain  $\mathbb{S}_h^N$  because  $\mathbb{S}_c^{N\perp} \cap \mathbb{S}_h^N = \mathbf{0}$ . Revising the argument of this inner-product form (1154), we get another flavor

$$\mathbf{D}(-V_N^T DV_N) = \begin{bmatrix} 0 \\ \delta(-V_N^T DV_N) \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & \delta(-V_N^T DV_N)^T \end{bmatrix} - 2 \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & -V_N^T DV_N \end{bmatrix} \quad (1157)$$

and we obtain mutual inversion of operators  $\mathbf{V}_N$  and  $\mathbf{D}$ , for  $D \in \mathbb{S}_h^N$

$$D = \mathbf{D}(\mathbf{V}_N(D)) \quad (1158)$$

$$-V_N^T DV_N = \mathbf{V}_N(\mathbf{D}(-V_N^T DV_N)) \quad (1159)$$

<sup>5.31</sup>Surjectivity of  $\mathbf{V}_N(D)$  is demonstrated via the Gram-form EDM operator  $\mathbf{D}(G)$ : Since  $\mathbb{S}_h^N = \mathbf{D}(\mathbb{S}_c^N)$  (1149), then for any  $Y \in \mathbb{S}^{N-1}$ ,  $-V_N^T \mathbf{D}(V_N^{\dagger T} Y V_N^{\dagger}/2) V_N = Y$ .

<sup>5.32</sup> $\mathcal{N}(\mathbf{V}_N) \supseteq \mathbb{S}_c^{N\perp}$  is apparent. There exists a linear mapping

$$T(\mathbf{V}_N(D)) \triangleq V_N^{\dagger T} \mathbf{V}_N(D) V_N^{\dagger} = -V D V \frac{1}{2} = \mathbf{V}(D)$$

such that

$$\mathcal{N}(T(\mathbf{V}_N)) = \mathcal{N}(\mathbf{V}) \supseteq \mathcal{N}(\mathbf{V}_N) \supseteq \mathbb{S}_c^{N\perp} = \mathcal{N}(\mathbf{V})$$

where the equality  $\mathbb{S}_c^{N\perp} = \mathcal{N}(\mathbf{V})$  is known (§E.7.2.0.2). ♦

or

$$\mathbb{S}_h^N = \mathbf{D}(\mathbf{V}_{\mathcal{N}}(\mathbb{S}_h^N)) \quad (1160)$$

$$-V_{\mathcal{N}}^T \mathbb{S}_h^N V_{\mathcal{N}} = \mathbf{V}_{\mathcal{N}}(\mathbf{D}(-V_{\mathcal{N}}^T \mathbb{S}_h^N V_{\mathcal{N}})) \quad (1161)$$

Substituting  $\Theta^T \Theta \leftarrow \Phi$  into inner-product form EDM definition (1103), any EDM may be expressed by the new flavor

$$\begin{aligned} \mathbf{D}(\Phi) &\triangleq \begin{bmatrix} 0 \\ \delta(\Phi) \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & \delta(\Phi)^T \end{bmatrix} - 2 \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \Phi \end{bmatrix} \in \mathbb{EDM}^N \\ &\Leftrightarrow \\ \Phi &\succeq 0 \end{aligned} \quad (1162)$$

where this  $\mathbf{D}$  is a linear surjective operator onto  $\mathbb{EDM}^N$  by definition, injective because it has no nullspace on domain  $\mathbb{S}_+^{N-1}$ . More broadly,  $\text{aff } \mathbf{D}(\mathbb{S}_+^{N-1}) = \mathbf{D}(\text{aff } \mathbb{S}_+^{N-1})$  (129),

$$\begin{aligned} \mathbb{S}_h^N &= \mathbf{D}(\mathbb{S}_+^{N-1}) \\ \mathbb{S}_+^{N-1} &= \mathbf{V}_{\mathcal{N}}(\mathbb{S}_h^N) \end{aligned} \quad (1163)$$

demonstrably isomorphisms, and by bijectivity of this inner-product form:

$$\mathbb{EDM}^N = \mathbf{D}(\mathbb{S}_+^{N-1}) \quad (1164)$$

$$\mathbb{S}_+^{N-1} = \mathbf{V}_{\mathcal{N}}(\mathbb{EDM}^N) \quad (1165)$$

## 5.7 Embedding in affine hull

The affine hull  $\mathcal{A}$  (79) of a point list  $\{x_\ell\}$  (arranged columnar in  $X \in \mathbb{R}^{n \times N}$  (77)) is identical to the affine hull of that polyhedron  $\mathcal{P}$  (87) formed from all convex combinations of the  $x_\ell$ ; [65, §2] [343, §17]

$$\mathcal{A} = \text{aff } X = \text{aff } \mathcal{P} \quad (1166)$$

Comparing hull definitions (79) and (87), it becomes obvious that the  $x_\ell$  and their convex hull  $\mathcal{P}$  are embedded in their unique affine hull  $\mathcal{A}$ ;

$$\mathcal{A} \supseteq \mathcal{P} \supseteq \{x_\ell\} \quad (1167)$$

Recall: *affine dimension r* is a lower bound on embedding, equal to dimension of the subspace parallel to that nonempty affine set  $\mathcal{A}$  in which the points are embedded. (§2.3.1) We define dimension of the convex hull  $\mathcal{P}$  to be the same as dimension  $r$  of the affine hull  $\mathcal{A}$  [343, §2], but  $r$  is not necessarily equal to rank of  $X$  (1186).

For the particular example illustrated in Figure 140,  $\mathcal{P}$  is the triangle in union with its relative interior while its three vertices constitute the entire list  $X$ . Affine hull  $\mathcal{A}$  is the unique plane that contains the triangle, so affine dimension  $r=2$  in that example while rank of  $X$  is 3. Were there only two points in Figure 140, then the affine hull would instead be the unique line passing through them;  $r$  would become 1 while rank would then be 2.

### 5.7.1 Determining affine dimension

Knowledge of affine dimension  $r$  becomes important because we lose any absolute offset common to all the generating  $x_\ell$  in  $\mathbb{R}^n$  when reconstructing convex polyhedra given only distance information. (§5.5.1) To calculate  $r$ , we first remove any offset that serves to

increase dimensionality of the subspace required to contain polyhedron  $\mathcal{P}$ ; subtracting any  $\alpha \in \mathcal{A}$  in the affine hull from every list member will work,

$$X - \alpha \mathbf{1}^T \quad (1168)$$

translating  $\mathcal{A}$  to the origin.<sup>5.33</sup>

$$\mathcal{A} - \alpha = \text{aff}(X - \alpha \mathbf{1}^T) = \text{aff}(X) - \alpha \quad (1169)$$

$$\mathcal{P} - \alpha = \text{conv}(X - \alpha \mathbf{1}^T) = \text{conv}(X) - \alpha \quad (1170)$$

Because (1166) and (1167) translate,

$$\mathbb{R}^n \supseteq \mathcal{A} - \alpha = \text{aff}(X - \alpha \mathbf{1}^T) = \text{aff}(\mathcal{P} - \alpha) \supseteq \mathcal{P} - \alpha \supseteq \{x_\ell - \alpha\} \quad (1171)$$

where from the previous relations it is easily shown

$$\text{aff}(\mathcal{P} - \alpha) = \text{aff}(\mathcal{P}) - \alpha \quad (1172)$$

Translating  $\mathcal{A}$  neither changes its dimension or the dimension of the embedded polyhedron  $\mathcal{P}$ ; (78)

$$r \triangleq \dim \mathcal{A} = \dim(\mathcal{A} - \alpha) \triangleq \dim(\mathcal{P} - \alpha) = \dim \mathcal{P} \quad (1173)$$

For any  $\alpha \in \mathbb{R}^n$ , (1169)-(1173) remain true. [343, p.4, p.12] Yet when  $\alpha \in \mathcal{A}$ , the affine set  $\mathcal{A} - \alpha$  becomes a unique subspace of  $\mathbb{R}^n$  in which the  $\{x_\ell - \alpha\}$  and their convex hull  $\mathcal{P} - \alpha$  are embedded (1171), and whose dimension is more easily calculated.

#### 5.7.1.0.1 Example. Translating first list-member to origin.

Subtracting the first member  $\alpha \triangleq x_1$  from every list member will translate their affine hull  $\mathcal{A}$  and their convex hull  $\mathcal{P}$  and, in particular,  $x_1 \in \mathcal{P} \subseteq \mathcal{A}$  to the origin in  $\mathbb{R}^n$ ; *videlicet*,

$$X - x_1 \mathbf{1}^T = X - X e_1 \mathbf{1}^T = X(I - e_1 \mathbf{1}^T) = X \begin{bmatrix} \mathbf{0} & \sqrt{2} V_N \end{bmatrix} \in \mathbb{R}^{n \times N} \quad (1174)$$

where  $V_N$  is defined in (1039), and  $e_1$  in (1049). Applying (1171) to (1174),

$$\mathbb{R}^n \supseteq \mathcal{R}(X V_N) = \mathcal{A} - x_1 = \text{aff}(X - x_1 \mathbf{1}^T) = \text{aff}(\mathcal{P} - x_1) \supseteq \mathcal{P} - x_1 \ni \mathbf{0} \quad (1175)$$

where  $X V_N \in \mathbb{R}^{n \times N-1}$ . Hence

$$r = \dim \mathcal{R}(X V_N) \quad (1176)$$

□

Since shifting the geometric center to the origin (§5.5.1.0.1) translates the affine hull to the origin as well, then it must also be true

$$r = \dim \mathcal{R}(X V) \quad (1177)$$

For any matrix whose range is  $\mathcal{R}(V) = \mathcal{N}(\mathbf{1}^T)$  we get the same result; *e.g.*,

$$r = \dim \mathcal{R}(X V_N^{\dagger T}) \quad (1178)$$

because

$$\mathcal{R}(X V) = \{Xz \mid z \in \mathcal{N}(\mathbf{1}^T)\} \quad (1179)$$

and  $\mathcal{R}(V) = \mathcal{R}(V_N) = \mathcal{R}(V_N^{\dagger T})$  (§E). These auxiliary matrices (§B.4.2) are more closely related;

$$V = V_N V_N^{\dagger} \quad (1815)$$

---

<sup>5.33</sup>Manipulation of hull functions  $\text{aff}$  and  $\text{conv}$  follows from their definitions.

### 5.7.1.1 Affine dimension $r$ versus rank

Now, suppose  $D$  is an EDM as defined by

$$\mathbf{D}(X) = \delta(X^T X) \mathbf{1}^T + \mathbf{1} \delta(X^T X)^T - 2X^T X \in \mathbb{EDM}^N \quad (1033)$$

and we premultiply by  $-V_N^T$  and postmultiply by  $V_N$ . Then because  $V_N^T \mathbf{1} = \mathbf{0}$  (1040), it is always true that

$$-V_N^T D V_N = 2V_N^T X^T X V_N = 2V_N^T G V_N \in \mathbb{S}^{N-1} \quad (1180)$$

where  $G$  is a Gram matrix. Similarly pre- and postmultiplying by  $V$  (confer (1054))

$$-V D V = 2V X^T X V = 2V G V \in \mathbb{S}^N \quad (1181)$$

always holds because  $V \mathbf{1} = \mathbf{0}$  (1805). Likewise, multiplying inner-product form EDM definition (1103), it always holds:

$$-V_N^T D V_N = \Theta^T \Theta \in \mathbb{S}^{N-1} \quad (1107)$$

For any matrix  $A$ ,  $\text{rank } A^T A = \text{rank } A = \text{rank } A^T$ . (1622) [228, §0.4]<sup>5.34</sup> So, by (1179), affine dimension

$$\begin{aligned} r &= \text{rank } X V = \text{rank } X V_N = \text{rank } X V_N^{\dagger T} = \text{rank } \Theta \\ &= \text{rank } V D V = \text{rank } V G V = \text{rank } V_N^T D V_N = \text{rank } V_N^T G V_N \end{aligned} \quad (1182)$$

By conservation of dimension, (§A.7.3.0.1)

$$r + \dim \mathcal{N}(V_N^T D V_N) = N-1 \quad (1183)$$

$$r + \dim \mathcal{N}(V D V) = N \quad (1184)$$

For  $D \in \mathbb{EDM}^N$

$$-V_N^T D V_N \succ 0 \Leftrightarrow r = N-1 \quad (1185)$$

but  $-V D V \not\succ 0$ . The general fact<sup>5.35</sup> (confer (1065))

$$r \leq \min\{n, N-1\} \quad (1186)$$

is evident from (1174) but can be visualized in the example illustrated in Figure 140. There we imagine a vector from the origin to each point in the list. Those three vectors are linearly independent in  $\mathbb{R}^3$ , but affine dimension  $r$  is 2 because the three points lie in a plane. When that plane is translated to the origin, it becomes the only subspace of dimension  $r=2$  that can contain the translated triangular polyhedron.

<sup>5.34</sup>For  $A \in \mathbb{R}^{m \times n}$ ,  $\mathcal{N}(A^T A) = \mathcal{N}(A)$ . [368, §3.3]

<sup>5.35</sup>  $\text{rank } X \leq \min\{n, N\}$

### 5.7.2 *Précis*

We collect expressions for affine dimension  $r$ : for list  $X \in \mathbb{R}^{n \times N}$  and Gram matrix  $G \in \mathbb{S}_+^N$

$$\begin{aligned} r &\triangleq \dim(\mathcal{P} - \alpha) = \dim \mathcal{P} = \dim \text{conv } X \quad (1173) \\ &= \dim(\mathcal{A} - \alpha) = \dim \mathcal{A} = \dim \text{aff } X \\ &= \text{rank}(X - x_1 \mathbf{1}^T) = \text{rank}(X - \frac{1}{N} X \mathbf{1} \mathbf{1}^T) \\ &= \text{rank } \Theta \quad (1105) \\ &= \text{rank } X V_{\mathcal{N}} = \text{rank } X V = \text{rank } X V_{\mathcal{N}}^{\dagger T} \\ &= \text{rank } X, \quad X e_1 = \mathbf{0} \quad \text{or} \quad X \mathbf{1} = \mathbf{0} \\ &= \text{rank } V_{\mathcal{N}}^T G V_{\mathcal{N}} = \text{rank } V G V = \text{rank } V_{\mathcal{N}}^{\dagger} G V_{\mathcal{N}} \\ &= \text{rank } G, \quad G e_1 = \mathbf{0} \quad (1050) \quad \text{or} \quad G \mathbf{1} = \mathbf{0} \quad (1054) \\ &= \text{rank } V_{\mathcal{N}}^T D V_{\mathcal{N}} = \text{rank } V D V = \text{rank } V_{\mathcal{N}}^{\dagger} D V_{\mathcal{N}} = \text{rank } V_{\mathcal{N}} (V_{\mathcal{N}}^T D V_{\mathcal{N}}) V_{\mathcal{N}}^T \\ &= \text{rank } \Lambda \quad (1273) \\ &= N - 1 - \dim \mathcal{N} \left( \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} \right) = \text{rank} \left[ \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} \right] - 2 \quad (1194) \end{aligned} \quad \left. \right\} D \in \mathbb{EDM}^N \quad (1187)$$

### 5.7.3 Eigenvalues of $-V D V$ versus $-V_{\mathcal{N}}^{\dagger} D V_{\mathcal{N}}$

Suppose for  $D \in \mathbb{EDM}^N$  we are given eigenvectors  $v_i \in \mathbb{R}^N$  of  $-V D V$  and corresponding eigenvalues  $\lambda \in \mathbb{R}^N$  so that

$$-V D V v_i = \lambda_i v_i, \quad i = 1 \dots N \quad (1188)$$

From these we can determine the eigenvectors and eigenvalues of  $-V_{\mathcal{N}}^{\dagger} D V_{\mathcal{N}}$ : Define

$$\nu_i \triangleq V_{\mathcal{N}}^{\dagger} v_i, \quad \lambda_i \neq 0 \quad (1189)$$

Then we have:

$$-V D V_{\mathcal{N}} V_{\mathcal{N}}^{\dagger} v_i = \lambda_i v_i \quad (1190)$$

$$-V_{\mathcal{N}}^{\dagger} V D V_{\mathcal{N}} \nu_i = \lambda_i V_{\mathcal{N}}^{\dagger} v_i \quad (1191)$$

$$-V_{\mathcal{N}}^{\dagger} D V_{\mathcal{N}} \nu_i = \lambda_i \nu_i \quad (1192)$$

the eigenvectors of  $-V_{\mathcal{N}}^{\dagger} D V_{\mathcal{N}}$  are given by (1189) while its corresponding nonzero eigenvalues are identical to those of  $-V D V$  although  $-V_{\mathcal{N}}^{\dagger} D V_{\mathcal{N}}$  is not necessarily positive semidefinite. In contrast,  $-V_{\mathcal{N}}^T D V_{\mathcal{N}}$  is positive semidefinite but its nonzero eigenvalues are generally different.

**5.7.3.0.1 Theorem.** *EDM rank versus affine dimension  $r$ .* [187, §3] [210, §3]  
[186, §3] For  $D \in \mathbb{EDM}^N$  (confer(1347))

$$1) \quad r = \text{rank}(D) - 1 \Leftrightarrow \mathbf{1}^T D^{\dagger} \mathbf{1} \neq 0$$

Points constituting a list  $X$  generating the polyhedron corresponding to  $D$  lie on the relative boundary of an  $r$ -dimensional *circumhypersphere* having

$$\begin{aligned} \text{diameter} &= \sqrt{2} (\mathbf{1}^T D^{\dagger} \mathbf{1})^{-1/2} \\ \text{circumcenter} &= \frac{X D^{\dagger} \mathbf{1}}{\mathbf{1}^T D^{\dagger} \mathbf{1}} \end{aligned} \quad (1193)$$

$$2) \quad r = \text{rank}(D) - 2 \Leftrightarrow \mathbf{1}^T D^{\dagger} \mathbf{1} = 0$$

There can be no circumhypersphere whose relative boundary contains a generating list for the corresponding polyhedron.

3) In *Cayley-Menger form* [126, §6.2] [96, §3.3] [54, §40] (§5.11.2),

$$r = N - 1 - \dim \mathcal{N} \left( \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} \right) = \text{rank} \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} - 2 \quad (1194)$$

Circumhyperspheres exist for  $r < \text{rank}(D) - 2$ . [386, §7]  $\diamond$

For all practical purposes, (1186)

$$\max\{0, \text{rank}(D) - 2\} \leq r \leq \min\{n, N - 1\} \quad (1195)$$

## 5.8 Euclidean metric *versus* matrix criteria

### 5.8.1 Nonnegativity property 1

When  $D = [d_{ij}]$  is an EDM (1033), then it is apparent from (1180)

$$2V_N^T X^T X V_N = -V_N^T D V_N \succeq 0 \quad (1196)$$

because for any matrix  $A$ ,  $A^T A \succeq 0$ .<sup>5.36</sup> We claim nonnegativity of the  $d_{ij}$  is enforced primarily by the matrix inequality (1196); *id est*,

$$\left. \begin{aligned} -V_N^T D V_N &\succeq 0 \\ D &\in \mathbb{S}_h^N \end{aligned} \right\} \Rightarrow d_{ij} \geq 0, \quad i \neq j \quad (1197)$$

(The matrix inequality to enforce strict positivity differs by a stroke of the pen. (1200))

We now support our claim: If any matrix  $A \in \mathbb{R}^{m \times m}$  is positive semidefinite, then its main diagonal  $\delta(A) \in \mathbb{R}^m$  must have all nonnegative entries. [181, §4.2]

Given  $D \in \mathbb{S}_h^N$

$$\begin{aligned} -V_N^T D V_N &= \\ \begin{bmatrix} d_{12} & \frac{1}{2}(d_{12} + d_{13} - d_{23}) & \frac{1}{2}(d_{1,i+1} + d_{1,j+1} - d_{i+1,j+1}) & \cdots & \frac{1}{2}(d_{12} + d_{1N} - d_{2N}) \\ \frac{1}{2}(d_{12} + d_{13} - d_{23}) & d_{13} & \frac{1}{2}(d_{1,i+1} + d_{1,j+1} - d_{i+1,j+1}) & \cdots & \frac{1}{2}(d_{13} + d_{1N} - d_{3N}) \\ \frac{1}{2}(d_{1,j+1} + d_{1,i+1} - d_{j+1,i+1}) & \frac{1}{2}(d_{1,j+1} + d_{1,i+1} - d_{j+1,i+1}) & d_{1,i+1} & \ddots & \frac{1}{2}(d_{14} + d_{1N} - d_{4N}) \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \frac{1}{2}(d_{12} + d_{1N} - d_{2N}) & \frac{1}{2}(d_{13} + d_{1N} - d_{3N}) & \frac{1}{2}(d_{14} + d_{1N} - d_{4N}) & \cdots & d_{1N} \end{bmatrix} \\ &= \frac{1}{2}(\mathbf{1} D_{1,2:N} + D_{2:N,1} \mathbf{1}^T - D_{2:N,2:N}) \in \mathbb{S}^{N-1} \end{aligned} \quad (1198)$$

where row, column indices  $i, j \in \{1 \dots N-1\}$ . [349] It follows:

$$\left. \begin{aligned} -V_N^T D V_N &\succeq 0 \\ D &\in \mathbb{S}_h^N \end{aligned} \right\} \Rightarrow \delta(-V_N^T D V_N) = \begin{bmatrix} d_{12} \\ d_{13} \\ \vdots \\ d_{1N} \end{bmatrix} \succeq 0 \quad (1199)$$

Multiplication of  $V_N$  by any permutation matrix  $\Xi$  has null effect on its range and nullspace. In other words, any permutation of the rows or columns of  $V_N$

<sup>5.36</sup>For  $A \in \mathbb{R}^{m \times n}$ ,  $A^T A \succeq 0 \Leftrightarrow \mathbf{y}^T A^T A \mathbf{y} = \|A\mathbf{y}\|^2 \geq 0$  for all  $\|\mathbf{y}\| = 1$ . When  $A$  is full-rank thin-or-square,  $A^T A \succ 0$ .

produces a basis for  $\mathcal{N}(\mathbf{1}^T)$ ; *id est*,  $\mathcal{R}(\Xi_r V_N) = \mathcal{R}(V_N \Xi_c) = \mathcal{R}(V_N) = \mathcal{N}(\mathbf{1}^T)$ . Hence,  $-V_N^T D V_N \succeq 0 \Leftrightarrow -V_N^T \Xi_r^T D \Xi_r V_N \succeq 0 \Leftrightarrow -\Xi_c^T V_N^T D V_N \Xi_c \succeq 0$ . Various permutation matrices will sift<sup>5.37</sup> remaining  $d_{ij}$  similarly to (1199) thereby proving their nonnegativity. Hence  $-V_N^T D V_N \succeq 0$  is a sufficient test for the first property (§5.2) of the Euclidean metric, nonnegativity. ♦

When affine dimension  $r$  equals 1, in particular, nonnegativity symmetry and hollowness become necessary and sufficient criteria satisfying matrix inequality (1196). (§6.5.0.0.1)

### 5.8.1.1 Strict positivity

Should we require the points in  $\mathbb{R}^n$  to be distinct, then entries of  $D$  off the main diagonal must be strictly positive  $\{d_{ij} > 0, i \neq j\}$  and only those entries along the main diagonal of  $D$  are 0. By similar argument, the strict matrix inequality is a sufficient test for strict positivity of Euclidean distance-square;

$$\left. \begin{array}{l} -V_N^T D V_N \succ 0 \\ D \in \mathbb{S}_h^N \end{array} \right\} \Rightarrow d_{ij} > 0, \quad i \neq j \quad (1200)$$

### 5.8.2 Triangle inequality property 4

In light of Kreyszig's observation [254, §1.1 prob.15] that properties 2 through 4 of the Euclidean metric (§5.2) together imply nonnegativity property 1,

$$2\sqrt{d_{jk}} = \sqrt{d_{jk}} + \sqrt{d_{kj}} \geq \sqrt{d_{jj}} = 0, \quad j \neq k \quad (1201)$$

nonnegativity criterion (1197) suggests that matrix inequality  $-V_N^T D V_N \succeq 0$  might somehow take on the role of triangle inequality; *id est*,

$$\left. \begin{array}{l} \delta(D) = \mathbf{0} \\ D^T = D \\ -V_N^T D V_N \succeq 0 \end{array} \right\} \Rightarrow \sqrt{d_{ij}} \leq \sqrt{d_{ik}} + \sqrt{d_{kj}}, \quad i \neq j \neq k \quad (1202)$$

We now show that is indeed the case: Let  $T$  be the *leading principal submatrix* in  $\mathbb{S}^2$  of  $-V_N^T D V_N$  (upper left  $2 \times 2$  submatrix from (1198));

$$T \triangleq \begin{bmatrix} d_{12} & \frac{1}{2}(d_{12}+d_{13}-d_{23}) \\ \frac{1}{2}(d_{12}+d_{13}-d_{23}) & d_{13} \end{bmatrix} \quad (1203)$$

Submatrix  $T$  must be positive (semi)definite whenever  $-V_N^T D V_N$  is. (§A.3.1.0.4, §5.8.3)  
Now we have,

$$\begin{aligned} -V_N^T D V_N \succeq 0 &\Rightarrow T \succeq 0 \Leftrightarrow \lambda_1 \geq \lambda_2 \geq 0 \\ -V_N^T D V_N \succ 0 &\Rightarrow T \succ 0 \Leftrightarrow \lambda_1 > \lambda_2 > 0 \end{aligned} \quad (1204)$$

where  $\lambda_1$  and  $\lambda_2$  are the eigenvalues of  $T$ , real due only to symmetry of  $T$ :

$$\begin{aligned} \lambda_1 &= \frac{1}{2} \left( d_{12} + d_{13} + \sqrt{d_{23}^2 - 2(d_{12} + d_{13})d_{23} + 2(d_{12}^2 + d_{13}^2)} \right) \in \mathbb{R} \\ \lambda_2 &= \frac{1}{2} \left( d_{12} + d_{13} - \sqrt{d_{23}^2 - 2(d_{12} + d_{13})d_{23} + 2(d_{12}^2 + d_{13}^2)} \right) \in \mathbb{R} \end{aligned} \quad (1205)$$

Nonnegativity of eigenvalue  $\lambda_1$  is guaranteed by only nonnegativity of the  $d_{ij}$  which in turn is guaranteed by matrix inequality (1197). Inequality between the eigenvalues in (1204) follows from only realness of the  $d_{ij}$ . Since  $\lambda_1$  always equals or exceeds  $\lambda_2$ ,

<sup>5.37</sup>Rule of thumb: If  $\Xi_r(i, 1) = 1$ , then  $\delta(-V_N^T \Xi_r^T D \Xi_r V_N) \in \mathbb{R}^{N-1}$  is some permutation of the  $i^{\text{th}}$  row or column of  $D$  excepting the 0 entry from the main diagonal.

conditions for positive (semi)definiteness of submatrix  $T$  can be completely determined by examining  $\lambda_2$  the smaller of its two eigenvalues. A triangle inequality is made apparent when we express  $T$  eigenvalue nonnegativity in terms of  $D$  matrix entries; *videlicet*,

$$\begin{aligned} T \succeq 0 &\Leftrightarrow \det T = \lambda_1 \lambda_2 \geq 0, \quad d_{12}, d_{13} \geq 0 \quad (\text{c}) \\ &\Leftrightarrow \lambda_2 \geq 0 \quad (\text{b}) \\ &\Leftrightarrow |\sqrt{d_{12}} - \sqrt{d_{23}}| \leq \sqrt{d_{13}} \leq \sqrt{d_{12}} + \sqrt{d_{23}} \quad (\text{a}) \end{aligned} \quad (1206)$$

Triangle inequality (1206a) (*confer*(1101)(1218)), in terms of three rooted entries from  $D$ , is equivalent to metric property 4

$$\begin{aligned} \sqrt{d_{13}} &\leq \sqrt{d_{12}} + \sqrt{d_{23}} \\ \sqrt{d_{23}} &\leq \sqrt{d_{12}} + \sqrt{d_{13}} \\ \sqrt{d_{12}} &\leq \sqrt{d_{13}} + \sqrt{d_{23}} \end{aligned} \quad (1207)$$

for the corresponding points  $x_1, x_2, x_3$  from some length- $N$  list.<sup>5.38</sup>

### 5.8.2.1 Comment

Given  $D$  whose dimension  $N$  equals or exceeds 3, there are  $N!/(3!(N-3)!)$  distinct triangle inequalities in total like (1101) that must be satisfied, of which each  $d_{ij}$  is involved in  $N-2$ , and each point  $x_i$  is in  $(N-1)!/(2!(N-1-2)!)$ . We have so far revealed only one of those triangle inequalities; namely, (1206a) that came from  $T$  (1203). Yet we claim if  $-V_N^T D V_N \succeq 0$  then all triangle inequalities will be satisfied simultaneously;

$$|\sqrt{d_{ik}} - \sqrt{d_{kj}}| \leq \sqrt{d_{ij}} \leq \sqrt{d_{ik}} + \sqrt{d_{kj}}, \quad i < k < j \quad (1208)$$

(There are no more.) To verify our claim, we must prove the matrix inequality  $-V_N^T D V_N \succeq 0$  to be a sufficient test of all the triangle inequalities; more efficient, we mention, for larger  $N$ :

**5.8.2.1.1 Shore.** The columns of  $\Xi_r V_N \Xi_c$  hold a basis for  $\mathcal{N}(\mathbf{1}^T)$  when  $\Xi_r$  and  $\Xi_c$  are permutation matrices. In other words, any permutation of the rows or columns of  $V_N$  leaves its range and nullspace unchanged; *id est*,  $\mathcal{R}(\Xi_r V_N \Xi_c) = \mathcal{R}(V_N) = \mathcal{N}(\mathbf{1}^T)$  (1040). Hence, two distinct matrix inequalities can be equivalent tests of the positive semidefiniteness of  $D$  on  $\mathcal{R}(V_N)$ ; *id est*,  $-V_N^T D V_N \succeq 0 \Leftrightarrow -(\Xi_r V_N \Xi_c)^T D (\Xi_r V_N \Xi_c) \succeq 0$ . By properly choosing permutation matrices,<sup>5.39</sup> the leading principal submatrix  $T_{\Xi} \in \mathbb{S}^2$  of  $-(\Xi_r V_N \Xi_c)^T D (\Xi_r V_N \Xi_c)$  may be loaded with the entries of  $D$  needed to test any particular triangle inequality (similarly to (1198)-(1206)). Because all the triangle inequalities can be individually tested using a test equivalent to the lone matrix inequality  $-V_N^T D V_N \succeq 0$ , it logically follows that the lone matrix inequality tests all those triangle inequalities simultaneously. We conclude that  $-V_N^T D V_N \succeq 0$  is a sufficient test for the fourth property of the Euclidean metric, triangle inequality. ♦

<sup>5.38</sup> Accounting for symmetry property 3, the fourth metric property demands three inequalities be satisfied per one of type (1206a). The first of those inequalities in (1207) is self evident from (1206a), while the two remaining follow from the left-hand side of (1206a) and the fact for scalars,  $|a| \leq b \Leftrightarrow a \leq b$  and  $-a \leq b$ .

<sup>5.39</sup> To individually test triangle inequality  $|\sqrt{d_{ik}} - \sqrt{d_{kj}}| \leq \sqrt{d_{ij}} \leq \sqrt{d_{ik}} + \sqrt{d_{kj}}$  for particular  $i, k, j$ , set  $\Xi_r(i, 1) = \Xi_r(k, 2) = \Xi_r(j, 3) = 1$  and  $\Xi_c = I$ .

### 5.8.2.2 Strict triangle inequality

Without exception, all the inequalities in (1206) and (1207) can be made strict while their corresponding implications remain true. The then strict inequality (1206a) or (1207) may be interpreted as a *strict triangle inequality* under which collinear arrangement of points is not allowed. [250, §24/6, p.322] Hence by similar reasoning,  $-V_N^T DV_N \succ 0$  is a sufficient test of all the strict triangle inequalities; *id est*,

$$\left. \begin{array}{l} \delta(D) = \mathbf{0} \\ D^T = D \\ -V_N^T DV_N \succ 0 \end{array} \right\} \Rightarrow \sqrt{d_{ij}} < \sqrt{d_{ik}} + \sqrt{d_{kj}}, \quad i \neq j \neq k \quad (1209)$$

### 5.8.3 $-V_N^T DV_N$ nesting

From (1203) observe that  $T = -V_N^T DV_N|_{N \leftarrow 3}$ . In fact, for  $D \in \mathbb{EDM}^N$ , the leading principal submatrices of  $-V_N^T DV_N$  form a nested sequence (by inclusion) whose members are individually positive semidefinite [181] [228] [368] and have the same form as  $T$ ; *videlicet*<sup>5.40</sup>,

$$-V_N^T DV_N|_{N \leftarrow 1} = [\emptyset] \quad (\text{o})$$

$$-V_N^T DV_N|_{N \leftarrow 2} = [d_{12}] \in \mathbb{S}_+ \quad (\text{a})$$

$$-V_N^T DV_N|_{N \leftarrow 3} = \begin{bmatrix} d_{12} & \frac{1}{2}(d_{12} + d_{13} - d_{23}) \\ \frac{1}{2}(d_{12} + d_{13} - d_{23}) & d_{13} \end{bmatrix} = T \in \mathbb{S}_+^2 \quad (\text{b})$$

$$-V_N^T DV_N|_{N \leftarrow 4} = \begin{bmatrix} d_{12} & \frac{1}{2}(d_{12} + d_{13} - d_{23}) & \frac{1}{2}(d_{12} + d_{14} - d_{24}) \\ \frac{1}{2}(d_{12} + d_{13} - d_{23}) & d_{13} & \frac{1}{2}(d_{13} + d_{14} - d_{34}) \\ \frac{1}{2}(d_{12} + d_{14} - d_{24}) & \frac{1}{2}(d_{13} + d_{14} - d_{34}) & d_{14} \end{bmatrix} \quad (\text{c})$$

⋮

$$-V_N^T DV_N|_{N \leftarrow i} = \begin{bmatrix} -V_N^T DV_N|_{N \leftarrow i-1} & \nu(i) \\ \nu(i)^T & d_{1i} \end{bmatrix} \in \mathbb{S}_+^{i-1} \quad (\text{d})$$

⋮

$$-V_N^T DV_N = \begin{bmatrix} -V_N^T DV_N|_{N \leftarrow N-1} & \nu(N) \\ \nu(N)^T & d_{1N} \end{bmatrix} \in \mathbb{S}_+^{N-1} \quad (\text{e}) \quad (1210)$$

where

$$\nu(i) \triangleq \frac{1}{2} \begin{bmatrix} d_{12} + d_{1i} - d_{2i} \\ d_{13} + d_{1i} - d_{3i} \\ \vdots \\ d_{1,i-1} + d_{1i} - d_{i-1,i} \end{bmatrix} \in \mathbb{R}^{i-2}, \quad i > 2 \quad (1211)$$

Hence, the leading principal submatrices of EDM  $D$  must also be EDMs<sup>5.41</sup>

<sup>5.40</sup>  $-V_N^T DV_N|_{N \leftarrow 1} = 0 \in \mathbb{S}_+^0$  (§B.4.1)

<sup>5.41</sup> In fact, each and every principal submatrix of an EDM  $D$  is another EDM. [264, §4.1]

Bordered symmetric matrices in the form (1210d) are known to have *intertwined* [368, §6.4] (or *interlaced* [228, §4.3] [364, §IV.4.1]) eigenvalues; (*confer* §5.11.1) that means, for the particular submatrices (1210a) and (1210b),

$$\lambda_2 \leq d_{12} \leq \lambda_1 \quad (1212)$$

where  $d_{12}$  is the eigenvalue of submatrix (1210a) and  $\lambda_1, \lambda_2$  are the eigenvalues of  $T$  (1210b) (1203). Intertwining in (1212) predicts that should  $d_{12}$  become 0, then  $\lambda_2$  must go to 0.<sup>5.42</sup> Eigenvalues are similarly intertwined for submatrices (1210b) and (1210c);

$$\gamma_3 \leq \lambda_2 \leq \gamma_2 \leq \lambda_1 \leq \gamma_1 \quad (1213)$$

where  $\gamma_1, \gamma_2, \gamma_3$  are the eigenvalues of submatrix (1210c). Intertwining likewise predicts that should  $\lambda_2$  become 0 (a possibility revealed in §5.8.3.1), then  $\gamma_3$  must go to 0. Combining results so far for  $N=2, 3, 4$ : (1212) (1213)

$$\gamma_3 \leq \lambda_2 \leq d_{12} \leq \lambda_1 \leq \gamma_1 \quad (1214)$$

The preceding logic extends by induction through the remaining members of the sequence (1210).

### 5.8.3.1 Tightening the triangle inequality

Now we apply Schur complement from §A.4 to tighten the triangle inequality from (1202) in case: cardinality  $N=4$ . We find that the gains by doing so are modest. From (1210) we identify:

$$\begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \triangleq -V_N^T D V_N|_{N \leftarrow 4} \quad (1215)$$

$$A \triangleq T = -V_N^T D V_N|_{N \leftarrow 3} \quad (1216)$$

both positive semidefinite by assumption, where  $B=\nu(4)$  (1211), and  $C=d_{14}$ . Using nonstrict  $CC^\dagger$ -form (1663),  $C \succeq 0$  by assumption (§5.8.1) and  $CC^\dagger=I$ . So by the *positive semidefinite ordering of eigenvalues theorem* (§A.3.1.0.1),

$$-V_N^T D V_N|_{N \leftarrow 4} \succeq 0 \Leftrightarrow T \succeq d_{14}^{-1} \nu(4) \nu(4)^T \Rightarrow \begin{cases} \lambda_1 \geq d_{14}^{-1} \|\nu(4)\|^2 \\ \lambda_2 \geq 0 \end{cases} \quad (1217)$$

where  $\{d_{14}^{-1} \|\nu(4)\|^2, 0\}$  are the eigenvalues of  $d_{14}^{-1} \nu(4) \nu(4)^T$  while  $\lambda_1, \lambda_2$  are the eigenvalues of  $T$ .

#### 5.8.3.1.1 Example. Small completion problem, II.

Applying the inequality for  $\lambda_1$  in (1217) to the *small completion problem* on page 343 Figure 141, the lower bound on  $\sqrt{d_{14}}$  (1.236 in (1026)) is tightened to 1.289. The correct value of  $\sqrt{d_{14}}$  to three significant figures is 1.414.  $\square$

---

<sup>5.42</sup>If  $d_{12}$  were 0, eigenvalue  $\lambda_2$  becomes 0 (1205) because  $d_{13}$  must then be equal to  $d_{23}$ ; *id est*,  $d_{12} = 0 \Leftrightarrow x_1 = x_2$ . (§5.4)

### 5.8.4 Affine dimension reduction in two dimensions

(confer §5.14.4) The leading principal  $2 \times 2$  submatrix  $T$  of  $-V_N^T D V_N$  has largest eigenvalue  $\lambda_1$  (1205) which is a convex function of  $D$ .<sup>5.43</sup>  $\lambda_1$  can never be 0 unless  $d_{12} = d_{13} = d_{23} = 0$ . Eigenvalue  $\lambda_1$  can never be negative while the  $d_{ij}$  are nonnegative. The remaining eigenvalue  $\lambda_2$  (1205) is a concave function of  $D$  that becomes 0 only at the upper and lower bounds of triangle inequality (1206a) and its equivalent forms: (confer(1208))

$$\begin{aligned} |\sqrt{d_{12}} - \sqrt{d_{23}}| &\leq \sqrt{d_{13}} \leq \sqrt{d_{12}} + \sqrt{d_{23}} & \text{(a)} \\ \Leftrightarrow |\sqrt{d_{12}} - \sqrt{d_{13}}| &\leq \sqrt{d_{23}} \leq \sqrt{d_{12}} + \sqrt{d_{13}} & \text{(b)} \\ \Leftrightarrow |\sqrt{d_{13}} - \sqrt{d_{23}}| &\leq \sqrt{d_{12}} \leq \sqrt{d_{13}} + \sqrt{d_{23}} & \text{(c)} \end{aligned} \quad (1218)$$

In between those bounds,  $\lambda_2$  is strictly positive; otherwise, it would be negative but prevented by the condition  $T \succeq 0$ .

When  $\lambda_2$  becomes 0, it means triangle  $\triangle_{123}$  has collapsed to a line segment; a potential reduction in affine dimension  $r$ . The same logic is valid for any particular principal  $2 \times 2$  submatrix of  $-V_N^T D V_N$ , hence applicable to other triangles.

## 5.9 Bridge: Convex polyhedra to EDMs

The criteria for the existence of an EDM include, by definition (1033) (1103), the properties imposed upon its entries  $d_{ij}$  by the Euclidean metric. From §5.8.1 and §5.8.2, we know there is a relationship of matrix criteria to those properties. Here is a snapshot of what we are sure: for  $i, j, k \in \{1 \dots N\}$  (confer §5.2)

$$\begin{array}{ll} \sqrt{d_{ij}} \geq 0, \quad i \neq j \\ \sqrt{d_{ij}} = 0, \quad i = j \\ \sqrt{d_{ij}} = \sqrt{d_{ji}} \\ \sqrt{d_{ij}} \leq \sqrt{d_{ik}} + \sqrt{d_{kj}}, \quad i \neq j \neq k \end{array} \Leftarrow \begin{array}{l} -V_N^T D V_N \succeq 0 \\ \delta(D) = \mathbf{0} \\ D^T = D \end{array} \quad (1219)$$

all implied by  $D \in \mathbb{EDM}^N$ . In words, these four Euclidean metric properties are necessary conditions for  $D$  to be a distance matrix. At the moment, we have no converse. As of concern in §5.3, we have yet to establish metric requirements beyond the four Euclidean metric properties that would allow  $D$  to be certified an EDM or might facilitate polyhedron or list reconstruction from an incomplete EDM. We deal with this problem in §5.14. Our present goal is to establish *ab initio* the necessary and sufficient matrix criteria that will subsume all the Euclidean metric properties and any further requirements<sup>5.44</sup> for all  $N > 1$  (§5.8.3); *id est*,

$$\left. \begin{array}{l} -V_N^T D V_N \succeq 0 \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1052)$$

<sup>5.43</sup>The largest eigenvalue of any symmetric matrix is always a convex function of its entries, while the smallest eigenvalue is always concave. [65, exmp.3.10] In our particular case, say  $\underline{d} \triangleq \begin{bmatrix} d_{12} \\ d_{13} \\ d_{23} \end{bmatrix} \in \mathbb{R}^3$ . Then

the Hessian (1936)  $\nabla^2 \lambda_1(\underline{d}) \succeq 0$  certifies convexity whereas  $\nabla^2 \lambda_2(\underline{d}) \preceq 0$  certifies concavity. Each Hessian has rank 1. The respective gradients  $\nabla \lambda_1(\underline{d})$  and  $\nabla \lambda_2(\underline{d})$  are nowhere  $\mathbf{0}$  and can be uniquely defined.

<sup>5.44</sup>Schoenberg [349, (1)] first extolled matrix product  $-V_N^T D V_N$  (1198) (predicated on symmetry and selfdistance) in 1935, specifically incorporating  $V_N$ , albeit algebraically. He showed: nonnegativity  $-y^T V_N^T D V_N y \geq 0$ ,  $\forall y \in \mathbb{R}^{N-1}$ , is necessary and sufficient for  $D$  to be an EDM. Gower [186, §3] remarks how surprising it is that such a fundamental property of Euclidean geometry was obtained so late.

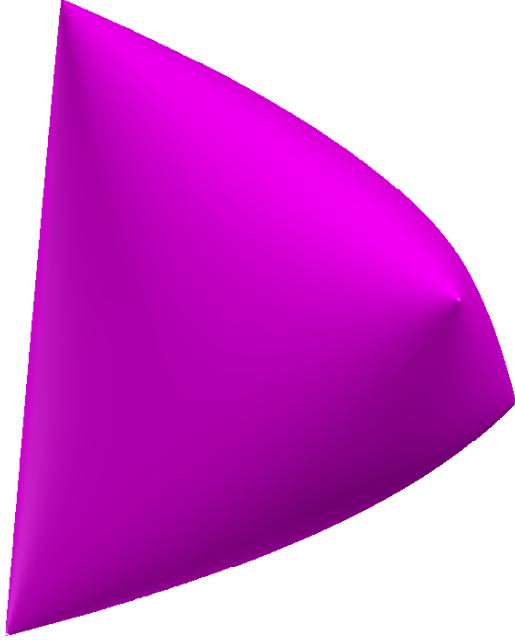


Figure 155: Elliptope  $\mathcal{E}^3$  in isometrically isomorphic  $\mathbb{R}^6$  (projected on  $\mathbb{R}^3$ ) is a convex body that appears to possess some kind of symmetry in this dimension; it resembles a malformed pillow in the shape of a bulging tetrahedron. Elliptope relative boundary is not *smooth* and comprises all set members (1221) having at least one 0 eigenvalue. [267, §2.1] This elliptope has an infinity of vertices, but there are only four vertices corresponding to a rank-1 matrix. Those  $yy^T$ , evident in the illustration, have binary vector  $y \in \mathbb{R}^3$  with entries in  $\{\pm 1\}$ .

or for EDM definition (1112),

$$\left. \begin{array}{l} \Omega \succeq 0 \\ \sqrt{\delta(d)} \succeq 0 \end{array} \right\} \Leftrightarrow D = \mathbf{D}(\Omega, d) \in \mathbb{EDM}^N \quad (1220)$$

### 5.9.1 Geometric arguments

**5.9.1.0.1 Definition.** *Elliptope:* [267] [264, §2.3] [126, §31.5] a unique bounded immutable convex Euclidean body in  $\mathbb{S}^n$ ; intersection of positive semidefinite cone  $\mathbb{S}_+^n$  with that set of  $n$  hyperplanes defined by unity main diagonal;

$$\mathcal{E}^n \triangleq \mathbb{S}_+^n \cap \{\Phi \in \mathbb{S}^n \mid \delta(\Phi) = 1\} \quad (1221)$$

a.k.a the set of all *correlation matrices* of dimension

$$\dim \mathcal{E}^n = n(n-1)/2 \text{ in } \mathbb{R}^{n(n+1)/2} \quad (1222)$$

An elliptope  $\mathcal{E}^n$  is not a polyhedron, in general, but has some polyhedral faces and an infinity of vertices.<sup>5.45</sup> Of those,  $2^{n-1}$  vertices (some extreme points of the elliptope) are

---

<sup>5.45</sup>Laurent defines vertex distinctly from the sense herein (§2.6.1.0.1); she defines *vertex* as a point with full-dimensional (nonempty interior) normal cone (§E.10.3.2.1). Her definition excludes point C in Figure 35, for example.

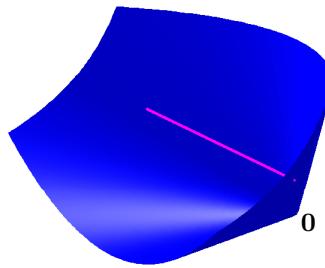


Figure 156: Ellipope  $\mathcal{E}^2$  in isometrically isomorphic  $\mathbb{R}^3$  is a line segment illustrated interior to positive semidefinite cone  $\mathbb{S}_+^2$  (Figure 46). Two vertices on boundary are rank-1 binary.

extreme directions  $yy^T$  of the positive semidefinite cone where entries of vector  $y \in \mathbb{R}^n$  belong to  $\{\pm 1\}$  and exercise every combination. Each of the remaining vertices has rank, greater than one, belonging to the set  $\{k > 0 \mid k(k+1)/2 \leq n\}$ . Each and every face of an ellipope is exposed.  $\triangle$

In fact, any positive semidefinite matrix whose entries belong to  $\{\pm 1\}$  is a rank-one correlation matrix; and *vice versa*.<sup>5.46</sup>

**5.9.1.0.2 Theorem.** *Ellipope vertices rank-one.* (confer §2.3.1.0.1) [143, §2.1.1]  
For  $Y \in \mathbb{S}^n$ ,  $y \in \mathbb{R}^n$ , and all  $i, j \in \{1 \dots n\}$

$$Y \succeq 0, \quad Y_{ij} \in \{\pm 1\} \quad \Leftrightarrow \quad Y = yy^T, \quad y_i \in \{\pm 1\} \quad (1223)$$

◊

The ellipope for dimension  $n=2$  is a line segment in isometrically isomorphic  $\mathbb{R}^{n(n+1)/2}$  (Figure 156). Obviously,  $\text{cone}(\mathcal{E}^n) \neq \mathbb{S}_+^n$ . The ellipope for dimension  $n=3$  is realized in Figure 155.

**5.9.1.0.3 Lemma.** *Hypersphere.* (confer bullet p.350) [18, §4]  
Matrix  $\Psi = [\Psi_{ij}] \in \mathbb{S}^N$  belongs to the ellipope in  $\mathbb{S}^N$  iff there exist  $N$  points  $p$  on the boundary of a hypersphere in  $\mathbb{R}^{\text{rank } \Psi}$  having radius 1 such that

$$\|p_i - p_j\|^2 = 2(1 - \Psi_{ij}), \quad i, j = 1 \dots N \quad (1224)$$

◊

There is a similar theorem for Euclidean distance matrices:

We derive matrix criteria for  $D$  to be an EDM, validating (1052) using simple geometry; distance to the polyhedron formed by the convex hull of a list of points (77) in Euclidean space  $\mathbb{R}^n$ .

---

<sup>5.46</sup>As there are few equivalent conditions for rank constraints, this device is rather important for relaxing integer, combinatorial, or Boolean problems.

#### 5.9.1.0.4 EDM assertion.

$D$  is a Euclidean distance matrix if and only if  $D \in \mathbb{S}_h^N$  and distances-square from the origin

$$\{\|p(y)\|^2 = -y^T V_N^T D V_N y \mid y \in \mathcal{S} - \beta\} \quad (1225)$$

correspond to points  $p$  in some bounded convex polyhedron

$$\mathcal{P} - \alpha = \{p(y) \mid y \in \mathcal{S} - \beta\} \quad (1226)$$

having  $N$  or fewer vertices embedded in an  $r$ -dimensional subspace  $\mathcal{A} - \alpha$  of  $\mathbb{R}^n$ , where  $\alpha \in \mathcal{A} = \text{aff } \mathcal{P}$  and where domain of linear surjection  $p(y)$  is the unit simplex  $\mathcal{S} \subset \mathbb{R}_+^{N-1}$  shifted such that its vertex at the origin is translated to  $-\beta$  in  $\mathbb{R}^{N-1}$ . When  $\beta = 0$ , then  $\alpha = x_1$ .  $\diamond$

In terms of  $V_N$ , the unit simplex (296) in  $\mathbb{R}^{N-1}$  has an equivalent representation:

$$\mathcal{S} = \{s \in \mathbb{R}^{N-1} \mid \sqrt{2}V_N s \succeq -e_1\} \quad (1227)$$

where  $e_1$  is as in (1049). Incidental to the *EDM assertion*, shifting the unit-simplex domain in  $\mathbb{R}^{N-1}$  translates the polyhedron  $\mathcal{P}$  in  $\mathbb{R}^n$ . Indeed, there is a map from vertices of the unit simplex to members of the list generating  $\mathcal{P}$ ;

$$\begin{aligned} p & : \mathbb{R}^{N-1} \rightarrow \mathbb{R}^n \\ p \left( \begin{Bmatrix} -\beta \\ e_1 - \beta \\ e_2 - \beta \\ \vdots \\ e_{N-1} - \beta \end{Bmatrix} \right) &= \begin{Bmatrix} x_1 - \alpha \\ x_2 - \alpha \\ x_3 - \alpha \\ \vdots \\ x_N - \alpha \end{Bmatrix} \end{aligned} \quad (1228)$$

#### 5.9.1.0.5 Proof. EDM assertion.

( $\Rightarrow$ ) We demonstrate that if  $D$  is an EDM, then each distance-square  $\|p(y)\|^2$  described by (1225) corresponds to a point  $p$  in some embedded polyhedron  $\mathcal{P} - \alpha$ . Assume  $D$  is indeed an EDM; *id est*,  $D$  can be made from some list  $X$  of  $N$  unknown points in Euclidean space  $\mathbb{R}^n$ ;  $D = \mathbf{D}(X)$  for  $X \in \mathbb{R}^{n \times N}$  as in (1033). Since  $D$  is translation invariant (§5.5.1), we may shift the affine hull  $\mathcal{A}$  of those unknown points to the origin as in (1168). Then take any point  $p$  in their convex hull (87);

$$\mathcal{P} - \alpha = \{p = (X - Xb\mathbf{1}^T)a \mid a^T\mathbf{1} = 1, a \succeq 0\} \quad (1229)$$

where  $\alpha = Xb \in \mathcal{A} \Leftrightarrow b^T\mathbf{1} = 1$ . Solutions to  $a^T\mathbf{1} = 1$  are:<sup>5.47</sup>

$$a \in \left\{ e_1 + \sqrt{2}V_N s \mid s \in \mathbb{R}^{N-1} \right\} \quad (1230)$$

where  $e_1$  is as in (1049). Similarly,  $b = e_1 + \sqrt{2}V_N \beta$ .

$$\begin{aligned} \mathcal{P} - \alpha &= \{p = X(I - (e_1 + \sqrt{2}V_N \beta)\mathbf{1}^T)(e_1 + \sqrt{2}V_N s) \mid \sqrt{2}V_N s \succeq -e_1\} \\ &= \{p = X\sqrt{2}V_N(s - \beta) \mid \sqrt{2}V_N s \succeq -e_1\} \end{aligned} \quad (1231)$$

that describes the domain of  $p(s)$  as the unit simplex

$$\mathcal{S} = \{s \mid \sqrt{2}V_N s \succeq -e_1\} \subset \mathbb{R}_+^{N-1} \quad (1227)$$

<sup>5.47</sup>Since  $\mathcal{R}(V_N) = \mathcal{N}(\mathbf{1}^T)$  and  $\mathcal{N}(\mathbf{1}^T) \perp \mathcal{R}(\mathbf{1})$ , then over all  $s \in \mathbb{R}^{N-1}$ ,  $V_N s$  is a hyperplane through the origin orthogonal to  $\mathbf{1}$ . Thus the solutions  $\{a\}$  constitute a hyperplane orthogonal to the vector  $\mathbf{1}$ , and offset from the origin in  $\mathbb{R}^N$  by any particular solution; in this case,  $a = e_1$ .

Making the substitution  $s - \beta \leftarrow y$

$$\mathcal{P} - \alpha = \{p = X\sqrt{2}V_N y \mid y \in \mathcal{S} - \beta\} \quad (1232)$$

Point  $p$  belongs to a convex polyhedron  $\mathcal{P} - \alpha$  embedded in an  $r$ -dimensional subspace of  $\mathbb{R}^n$  because the convex hull of any list forms a polyhedron, and because the translated affine hull  $\mathcal{A} - \alpha$  contains the translated polyhedron  $\mathcal{P} - \alpha$  (1171) and the origin (when  $\alpha \in \mathcal{A}$ ), and because  $\mathcal{A}$  has dimension  $r$  by definition (1173). Now, any distance-square from the origin to the polyhedron  $\mathcal{P} - \alpha$  can be formulated

$$\{p^T p = \|p\|^2 = 2y^T V_N^T X^T X V_N y \mid y \in \mathcal{S} - \beta\} \quad (1233)$$

Applying (1180) to (1233) we get (1225).

( $\Leftarrow$ ) To validate the *EDM assertion* in the reverse direction, we prove: If each distance-square  $\|p(y)\|^2$  (1225) on the shifted unit-simplex  $\mathcal{S} - \beta \subset \mathbb{R}^{N-1}$  corresponds to a point  $p(y)$  in some embedded polyhedron  $\mathcal{P} - \alpha$ , then  $D$  is an EDM. The  $r$ -dimensional subspace  $\mathcal{A} - \alpha \subseteq \mathbb{R}^n$  is spanned by

$$p(\mathcal{S} - \beta) = \mathcal{P} - \alpha \quad (1234)$$

because  $\mathcal{A} - \alpha = \text{aff}(\mathcal{P} - \alpha) \supseteq \mathcal{P} - \alpha$  (1171). So, outside domain  $\mathcal{S} - \beta$  of linear surjection  $p(y)$ , simplex complement  $\setminus \mathcal{S} - \beta \subset \mathbb{R}^{N-1}$  must contain domain of the distance-square  $\|p(y)\|^2 = p(y)^T p(y)$  to remaining points in subspace  $\mathcal{A} - \alpha$ ; *id est*, to the polyhedron's relative exterior  $\setminus \mathcal{P} - \alpha$ . For  $\|p(y)\|^2$  to be nonnegative on the entire subspace  $\mathcal{A} - \alpha$ ,  $-V_N^T D V_N$  must be positive semidefinite and is assumed symmetric;<sup>5.48</sup>

$$-V_N^T D V_N \triangleq \Theta_p^T \Theta_p \quad (1235)$$

where<sup>5.49</sup>  $\Theta_p \in \mathbb{R}^{m \times N-1}$  for some  $m \geq r$ . Because  $p(\mathcal{S} - \beta)$  is a convex polyhedron, it is necessarily a set of linear combinations of points from some length- $N$  list because every convex polyhedron having  $N$  or fewer vertices can be generated that way (§2.12.2). Equivalent to (1225) are

$$\{p^T p \mid p \in \mathcal{P} - \alpha\} = \{p^T p = y^T \Theta_p^T \Theta_p y \mid y \in \mathcal{S} - \beta\} \quad (1236)$$

Because  $p \in \mathcal{P} - \alpha$  may be found by factoring (1236), the list  $\Theta_p$  is found by factoring (1235). A unique EDM can be made from that list using inner-product form definition  $\mathbf{D}(\Theta)|_{\Theta=\Theta_p}$  (1103). That EDM will be identical to  $D$  if  $\delta(D)=\mathbf{0}$ , by injectivity of  $\mathbf{D}$  (1154).  $\blacklozenge$

### 5.9.2 Necessity and sufficiency

From (1196) we learned that matrix inequality  $-V_N^T D V_N \succeq 0$  is a necessary test for  $D$  to be an EDM. In §5.9.1, the connection between convex polyhedra and EDMs was pronounced by the *EDM assertion*; the matrix inequality together with  $D \in \mathbb{S}_h^N$  became a sufficient test when the *EDM assertion* demanded that every bounded convex polyhedron have a corresponding EDM. For all  $N > 1$  (§5.8.3), the matrix criteria for the existence of an EDM in (1052), (1220), and (1028) are therefore necessary and sufficient and subsume all the Euclidean metric properties and further requirements.

### 5.9.3 Example revisited

Now we apply the necessary and sufficient EDM criteria (1052) to an earlier problem.

<sup>5.48</sup>The antisymmetric part  $(-V_N^T D V_N - (-V_N^T D V_N)^T)/2$  is annihilated by  $\|p(y)\|^2$ . By the same reasoning, any positive (semi)definite matrix  $A$  is generally assumed symmetric because only the symmetric part  $(A + A^T)/2$  survives the test  $y^T A y \geq 0$ . [228, §7.1]

<sup>5.49</sup> $A^T = A \succeq 0 \Leftrightarrow A = R^T R$  for some real matrix  $R$ . [368, §6.3]

**5.9.3.0.1 Example.** *Small completion problem, III.* (confer §5.8.3.1.1)

Continuing Example 5.3.0.0.2 pertaining to Figure 141 where  $N = 4$ , distance-square  $d_{14}$  is ascertainable from the matrix inequality  $-V_N^T D V_N \succeq 0$ . Because all distances in (1025) are known except  $\sqrt{d_{14}}$ , we may simply calculate the smallest eigenvalue of  $-V_N^T D V_N$  over a range of  $d_{14}$  as in Figure 157. We observe a unique value of  $d_{14}$  satisfying (1052) where the abscissa axis is tangent to the hypograph of the smallest eigenvalue. Since the smallest eigenvalue of a symmetric matrix is known to be a concave function (§5.8.4), we calculate its second partial derivative with respect to  $d_{14}$  evaluated at 2 and find  $-1/3$ . We conclude there are no other satisfying values of  $d_{14}$ . Further, that value of  $d_{14}$  does not meet an upper or lower bound of a triangle inequality like (1208), so neither does it cause collapse of any triangle. Because the smallest eigenvalue is 0, affine dimension  $r$  of any point list corresponding to  $D$  cannot exceed  $N - 2$ . (§5.7.1.1)  $\square$

## 5.10 EDM-entry composition

Laurent [264, §2.3] applies results from Schoenberg, 1938 [350], to show certain nonlinear compositions of individual EDM entries yield EDMs; in particular,

$$\begin{aligned} D \in \mathbb{EDM}^N &\Leftrightarrow [1 - e^{-\alpha d_{ij}}] \in \mathbb{EDM}^N \quad \forall \alpha > 0 \quad (\text{a}) \\ &\Leftrightarrow [e^{-\alpha d_{ij}}] \in \mathcal{E}^N \quad \forall \alpha > 0 \quad (\text{b}) \end{aligned} \quad (1237)$$

where  $D = [d_{ij}]$  and  $\mathcal{E}^N$  is the ellipope (1221).

**5.10.0.0.1 Proof.** (Monique Laurent, 2003)

[350] (confer [254])

**Lemma 2.1.** from *A Tour d'Horizon ... on Completion Problems.* [264]

For  $D = [d_{ij}, i, j = 1 \dots N] \in \mathbb{S}_h^N$  and  $\mathcal{E}^N$  the ellipope in  $\mathbb{S}^N$  (§5.9.1.0.1), the following assertions are equivalent:

- (i)  $D \in \mathbb{EDM}^N$
- (ii)  $e^{-\alpha D} \triangleq [e^{-\alpha d_{ij}}] \in \mathcal{E}^N$  for all  $\alpha > 0$
- (iii)  $\mathbf{1}\mathbf{1}^T - e^{-\alpha D} \triangleq [1 - e^{-\alpha d_{ij}}] \in \mathbb{EDM}^N$  for all  $\alpha > 0$

$\diamond$

1) Equivalence of Lemma 2.1 (i) (ii) is stated in Schoenberg's Theorem 1 [350, p.527].

2) (ii)  $\Rightarrow$  (iii) can be seen from the statement in the beginning of section 3, saying that a distance space embeds in  $L_2$  iff some associated matrix is PSD. We reformulate it:

Let  $d = (d_{ij})_{i,j=0,1\dots N}$  be a distance space on  $N+1$  points (*i.e.*, symmetric hollow matrix of order  $N+1$ ) and let  $p = (p_{ij})_{i,j=1\dots N}$  be the symmetric matrix of order  $N$  related by:

$$\begin{aligned} (\mathbf{A}) \quad 2p_{ij} &= d_{0i} + d_{0j} - d_{ij} \quad \text{for } i, j = 1 \dots N \\ &\quad \text{or equivalently} \end{aligned}$$

$$(\mathbf{B}) \quad d_{0i} = p_{ii}, \quad d_{ij} = p_{ii} + p_{jj} - 2p_{ij} \quad \text{for } i, j = 1 \dots N$$

Then  $d$  embeds in  $L_2$  iff  $p$  is a positive semidefinite matrix iff  $d$  is of negative type (second half page 525/top of page 526 in [350]).

For the implication from (ii) to (iii), set:  $p = e^{-\alpha d}$  and define  $d'$  from  $p$  using (B) above. Then  $d'$  is a distance space on  $N+1$  points that embeds in  $L_2$ .

Thus its subspace of  $N$  points also embeds in  $L_2$  and is precisely  $1 - e^{-\alpha d}$ .

Note that (iii)  $\Rightarrow$  (ii) cannot be read immediately from this argument since (iii) involves the subdistance of  $d'$  on  $N$  points (and not the full  $d'$  on  $N+1$  points).

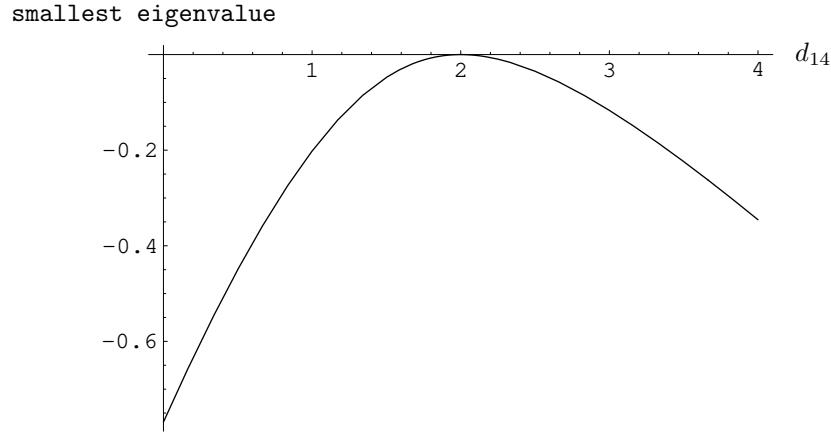


Figure 157: Smallest eigenvalue of  $-V_N^T D V_N$  makes it a PSD matrix for only one value of  $d_{14}$ : 2.

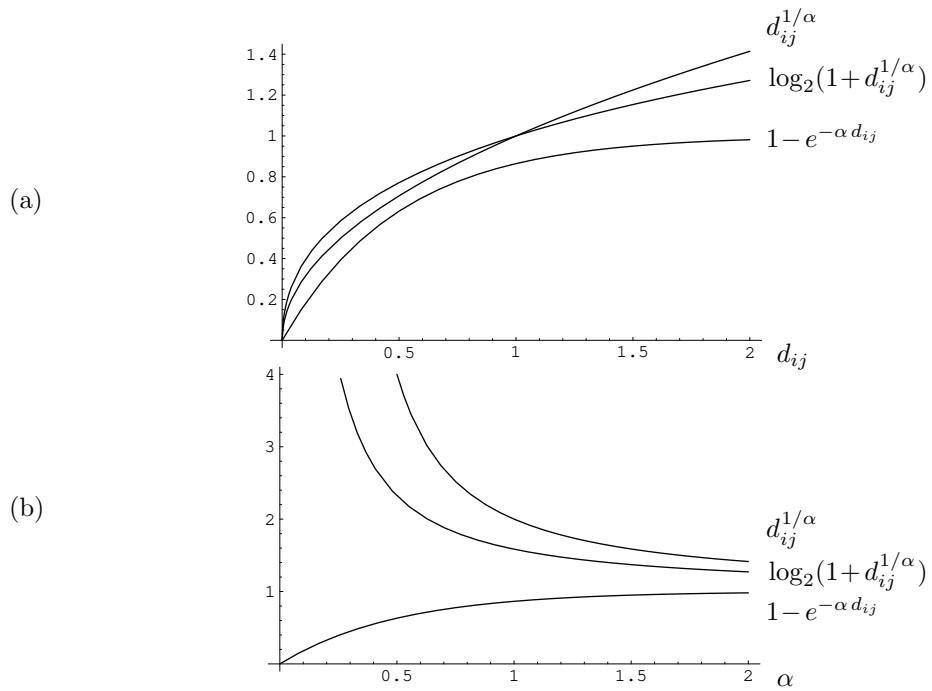


Figure 158: Some entrywise EDM compositions: (a)  $\alpha = 2$ . Concave nondecreasing in  $d_{ij}$ . (b) Trajectory convergence in  $\alpha$  for  $d_{ij} = 2$ .

- 3) Show (iii)  $\Rightarrow$  (i) by using the series expansion of the function  $1 - e^{-\alpha d}$ : the constant term cancels,  $\alpha$  factors out; there remains a summation of  $d$  plus a multiple of  $\alpha$ . Letting  $\alpha$  go to 0 gives the result.

This is not explicitly written in Schoenberg, but he also uses such an argument; expansion of the exponential function then  $\alpha \rightarrow 0$  (first proof on [350, p.526]). ♦

Schoenberg's results [350, §6 thm.5] (confer [254, p.108-109]) also suggest certain finite positive roots of EDM entries produce EDMs; specifically,

$$D \in \mathbb{EDM}^N \Leftrightarrow [d_{ij}^{1/\alpha}] \in \mathbb{EDM}^N \quad \forall \alpha > 1 \quad (1238)$$

The special case  $\alpha = 2$  is of interest because it corresponds to absolute distance; *e.g.*,

$$D \in \mathbb{EDM}^N \Rightarrow \sqrt[2]{D} \in \mathbb{EDM}^N \quad (1239)$$

Assuming that points constituting a corresponding list  $X$  are distinct (1200), then it follows: for  $D \in \mathbb{S}_h^N$

$$\lim_{\alpha \rightarrow \infty} [d_{ij}^{1/\alpha}] = \lim_{\alpha \rightarrow \infty} [1 - e^{-\alpha d_{ij}}] = -E \triangleq \mathbf{1}\mathbf{1}^T - I \quad (1240)$$

Negative elementary matrix  $-E$  (§B.3) is: relatively interior to the EDM cone (§6.5), on its axis, and terminal to respective trajectories (1237a) and (1238) as functions of  $\alpha$ . Both trajectories are confined to the EDM cone; in engineering terms, the EDM cone is an *invariant set* [346] to either trajectory. Further, if  $D$  is not an EDM but for some particular  $\alpha_p$  it becomes an EDM, then for all greater values of  $\alpha$  it remains an EDM. ▼

#### 5.10.0.0.2 Exercise. Concave nondecreasing EDM-entry composition.

Given EDM  $D = [d_{ij}]$ , empirical evidence suggests that the composition  $[\log_2(1 + d_{ij}^{1/\alpha})]$  is also an EDM for each fixed  $\alpha \geq 1$  [*sic*]. Its concavity in  $d_{ij}$  is illustrated in Figure 158 together with functions from (1237a) and (1238). Prove whether it holds more generally: Any concave nondecreasing composition of individual EDM entries  $d_{ij}$  on  $\mathbb{R}_+$  produces another EDM. ▼

#### 5.10.0.0.3 Exercise. Taxicab distance matrix as EDM.

Determine whether taxicab distance matrices ( $\mathbf{D}_1(X)$  in Example 3.10.0.0.2) are all numerically equivalent to EDMs. Explain why or why not. ▼

### 5.10.1 EDM by ellotope

(confer(1059)) For some  $\kappa \in \mathbb{R}_+$  and  $C \in \mathbb{S}_+^N$  in ellotope  $\mathcal{E}^N$  (§5.9.1.0.1), Alfakih asserts: any given EDM  $D$  is expressible [10] [126, §31.5]

$$D = \kappa(\mathbf{1}\mathbf{1}^T - C) \in \mathbb{EDM}^N \quad (1241)$$

This expression exhibits nonlinear combination of variables  $\kappa$  and  $C$ . We therefore propose a different expression requiring redefinition of the ellotope (1221) by scalar parametrization;

$$\mathcal{E}_t^n \triangleq \mathbb{S}_+^n \cap \{\Phi \in \mathbb{S}^n \mid \delta(\Phi) = t\mathbf{1}\} \quad (1242)$$

where, of course,  $\mathcal{E}^n = \mathcal{E}_1^n$ . Then any given EDM  $D$  is expressible

$$D = t\mathbf{1}\mathbf{1}^T - \mathfrak{C} \in \mathbb{EDM}^N \quad (1243)$$

which is linear in variables  $t \in \mathbb{R}_+$  and  $\mathfrak{C} \in \mathcal{E}_t^n$ .

## 5.11 EDM indefiniteness

By known result (§A.7.2) regarding a 0-valued entry on the main diagonal of a symmetric positive semidefinite matrix, there can be no positive or negative semidefinite EDM except the  $\mathbf{0}$  matrix because  $\text{EDM}^N \subseteq \mathbb{S}_h^N$  (1032) and

$$\mathbb{S}_h^N \cap \mathbb{S}_+^N = \mathbf{0} \quad (1244)$$

the origin. So when  $D \in \text{EDM}^N$ , there can be no factorization  $D = A^T A$  or  $-D = A^T A$ . [368, §6.3] Hence eigenvalues of an EDM are neither all nonnegative or all nonpositive; an EDM is indefinite and possibly invertible.

### 5.11.1 EDM eigenvalues, congruence transformation

For any symmetric  $-D$ , we can characterize its eigenvalues by congruence transformation: [368, §6.3]

$$-W^T D W = -\begin{bmatrix} V_N^T \\ \mathbf{1}^T \end{bmatrix} D \begin{bmatrix} V_N & \mathbf{1} \end{bmatrix} = -\begin{bmatrix} V_N^T D V_N & V_N^T D \mathbf{1} \\ \mathbf{1}^T D V_N & \mathbf{1}^T D \mathbf{1} \end{bmatrix} \in \mathbb{S}^N \quad (1245)$$

Because

$$W \triangleq [V_N \ \mathbf{1}] \in \mathbb{R}^{N \times N} \quad (1246)$$

is full-rank, then (1668)

$$\text{inertia}(-D) = \text{inertia}(-W^T D W) \quad (1247)$$

the congruence (1245) has the same number of positive, zero, and negative eigenvalues as  $-D$ . Further, if we denote by  $\{\gamma_i, i=1 \dots N-1\}$  the eigenvalues of  $-V_N^T D V_N$  and denote eigenvalues of the congruence  $-W^T D W$  by  $\{\zeta_i, i=1 \dots N\}$  and if we arrange each respective set of eigenvalues in nonincreasing order, then by theory of *interlacing eigenvalues for bordered symmetric matrices* [228, §4.3] [368, §6.4] [364, §IV.4.1]

$$\zeta_N \leq \gamma_{N-1} \leq \zeta_{N-1} \leq \gamma_{N-2} \leq \dots \leq \gamma_2 \leq \zeta_2 \leq \gamma_1 \leq \zeta_1 \quad (1248)$$

When  $D \in \text{EDM}^N$ , then  $\gamma_i \geq 0 \forall i$  (1604) because  $-V_N^T D V_N \succeq 0$  as we know. That means the congruence must have  $N-1$  nonnegative eigenvalues;  $\zeta_i \geq 0, i=1 \dots N-1$ . The remaining eigenvalue  $\zeta_N$  cannot be nonnegative because then  $-D$  would be positive semidefinite, an impossibility; so  $\zeta_N < 0$ . By congruence, nontrivial  $-D$  must therefore have exactly one negative eigenvalue;<sup>5.50</sup> [126, §2.4.5]

$$D \in \text{EDM}^N \Rightarrow \begin{cases} \lambda(-D)_i \geq 0, & i=1 \dots N-1 \\ \left( \sum_{i=1}^N \lambda(-D)_i = 0 \right) \\ D \in \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N} \end{cases} \quad (1249)$$

where the  $\lambda(-D)_i$  are nonincreasingly ordered eigenvalues of  $-D$  whose sum must be 0 only because  $\text{tr } D = 0$  [368, §5.1]. The eigenvalue summation condition, therefore, can be considered redundant. Even so, all these conditions are insufficient to determine whether some given  $H \in \mathbb{S}_h^N$  is an EDM; as shown by counterexample.<sup>5.51</sup>

<sup>5.50</sup>All entries of the corresponding eigenvector must have the same sign, with respect to each other, [97, p.116] because that eigenvector is the *Perron vector* corresponding to *spectral radius*; [228, §8.3.1] the predominant characteristic of square nonnegative matrices. Unlike positive semidefinite matrices, nonnegative matrices are guaranteed only to have at least one nonnegative eigenvalue.

<sup>5.51</sup>When  $N=3$ , for example, the symmetric hollow matrix

$$H = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 5 \\ 1 & 5 & 0 \end{bmatrix} \in \mathbb{S}_h^3 \cap \mathbb{R}_+^{3 \times 3}$$

**5.11.1.0.1 Exercise.** *Spectral inequality.*

Prove whether it holds: for  $D = [d_{ij}] \in \mathbb{EDM}^N$

$$\lambda(-D)_1 \geq d_{ij} \geq \lambda(-D)_{N-1} \quad \forall i \neq j \quad (1250)$$

▼

**5.11.1.0.2 Definition.** *Spectral cone  $\mathcal{K}_\lambda$ .*

A convex cone containing all *eigenspectra* corresponding to some given set of matrices is called a *spectral cone*. △

**5.11.1.0.3 Definition.** *Eigenspectrum.* [254, p.365] [364, p.26] (confer §A.5.0.1)

The eigenvalues of a matrix, including duplicates, are referred to as its *eigenspectrum*. △

Any positive semidefinite matrix, for example, possesses a vector (or nonincreasing list) of nonnegative eigenvalues corresponding to an eigenspectrum contained in a spectral cone  $\mathcal{K}_\lambda$  that is a nonnegative orthant (or monotone nonnegative cone).

**5.11.2 Spectral cones  $\mathcal{K}_\lambda$  for distance matrices**

Denoting the eigenvalues of Cayley-Menger matrix  $\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} \in \mathbb{S}^{N+1}$  by

$$\lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix}\right) \in \mathbb{R}^{N+1} \quad (1251)$$

we have the Cayley-Menger form (§5.7.3.0.1) of necessary and sufficient conditions for  $D \in \mathbb{EDM}^N$  from the literature: [210, §3]<sup>5.52</sup> [81, §3] [126, §6.2] (confer (1052) (1028))

$$D \in \mathbb{EDM}^N \Leftrightarrow \left\{ \begin{array}{l} \lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix}\right)_i \geq 0, \quad i = 1 \dots N \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow \left\{ \begin{array}{l} -V_N^T D V_N \succeq 0 \\ D \in \mathbb{S}_h^N \end{array} \right\} \quad (1252)$$

These conditions say the Cayley-Menger form has one and only one negative eigenvalue. When  $D$  is an EDM, eigenvalues  $\lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix}\right)$  belong to that particular orthant in  $\mathbb{R}^{N+1}$  having the  $N+1^{\text{th}}$  coordinate as sole negative coordinate:<sup>5.53</sup>

$$\begin{bmatrix} \mathbb{R}_+^N \\ \mathbb{R}_- \end{bmatrix} = \text{cone} \{e_1, e_2, \dots, e_N, -e_{N+1}\} \quad (1253)$$

**5.11.2.1 Cayley-Menger versus Schoenberg**

Connection to the Schoenberg criterion (1052) is made when the Cayley-Menger form is further partitioned:

$$\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} & \begin{bmatrix} \mathbf{1}^T \\ -D_{1,2:N} \end{bmatrix} \\ [\mathbf{1} \quad -D_{2:N,1}] & -D_{2:N,2:N} \end{bmatrix} \quad (1254)$$

is not an EDM, although  $\lambda(-H) = [5 \ 0.3723 \ -5.3723]^T$  conforms to (1249).

<sup>5.52</sup> Recall: for  $D \in \mathbb{S}_h^N$ ,  $-V_N^T D V_N \succeq 0$  subsumes nonnegativity property 1 (§5.8.1).

<sup>5.53</sup> Empirically, all except one entry of the corresponding eigenvector have the same sign with respect to each other.

Matrix  $D \in \mathbb{S}_h^N$  is an EDM if and only if the Schur complement of  $\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$  (§A.4) in this partition is positive semidefinite; [18, §1] [244, §3] *id est*, (confer(1198))

$$\begin{aligned} D \in \mathbb{EDM}^N \\ \Leftrightarrow \\ -D_{2:N, 2:N} - [\mathbf{1} \quad -D_{2:N, 1}] \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{1}^T \\ -D_{1, 2:N} \end{bmatrix} = -2V_{\mathcal{N}}^T DV_{\mathcal{N}} \succeq 0 \\ \text{and} \\ D \in \mathbb{S}_h^N \end{aligned} \quad (1255)$$

Positive semidefiniteness of that Schur complement insures nonnegativity ( $D \in \mathbb{R}_+^{N \times N}$ , §5.8.1), whereas *complementary inertia* (1670) insures existence of that lone negative eigenvalue of the Cayley-Menger form.

Now we apply results from chapter 2 with regard to polyhedral cones and their duals.

### 5.11.2.2 Ordered eigenspectra

Conditions (1252) specify eigenvalue membership to  $\mathcal{K}_{\lambda}$  the smallest pointed polyhedral *spectral cone* for  $\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbb{EDM}^N \end{bmatrix}$ :

$$\begin{aligned} \mathcal{K}_{\lambda} &\triangleq \{\zeta \in \mathbb{R}^{N+1} \mid \zeta_1 \geq \zeta_2 \geq \cdots \geq \zeta_N \geq 0 \geq \zeta_{N+1}, \mathbf{1}^T \zeta = 0\} \\ &= \mathcal{K}_{\mathcal{M}} \cap \left[ \begin{array}{c} \mathbb{R}_+^N \\ \mathbb{R}_- \end{array} \right] \cap \partial \mathcal{H} \\ &= \lambda \left( \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbb{EDM}^N \end{bmatrix} \right) \end{aligned} \quad (1256)$$

where

$$\partial \mathcal{H} \triangleq \{\zeta \in \mathbb{R}^{N+1} \mid \mathbf{1}^T \zeta = 0\} \quad (1257)$$

is a hyperplane through the origin, and  $\mathcal{K}_{\mathcal{M}}$  is the monotone cone (§2.13.10.4.3, implying ordered eigenspectra) which is full-dimensional but is not pointed;

$$\mathcal{K}_{\mathcal{M}} = \{\zeta \in \mathbb{R}^{N+1} \mid \zeta_1 \geq \zeta_2 \geq \cdots \geq \zeta_{N+1}\} \quad (441)$$

$$\mathcal{K}_{\mathcal{M}}^* = \{[e_1 - e_2 \quad e_2 - e_3 \quad \cdots \quad e_N - e_{N+1}] a \mid a \succeq 0\} \subset \mathbb{R}^{N+1} \quad (442)$$

So because of the hyperplane,

$$\dim \text{aff } \mathcal{K}_{\lambda} = \dim \partial \mathcal{H} = N \quad (1258)$$

indicating that spectral cone  $\mathcal{K}_{\lambda}$  is not full-dimensional. Defining

$$A \triangleq \begin{bmatrix} e_1^T - e_2^T \\ e_2^T - e_3^T \\ \vdots \\ e_N^T - e_{N+1}^T \end{bmatrix} \in \mathbb{R}^{N \times N+1}, \quad B \triangleq \begin{bmatrix} e_1^T \\ e_2^T \\ \vdots \\ e_N^T \\ -e_{N+1}^T \end{bmatrix} \in \mathbb{R}^{N+1 \times N+1} \quad (1259)$$

we have the halfspace-description:

$$\mathcal{K}_{\lambda} = \{\zeta \in \mathbb{R}^{N+1} \mid A \zeta \succeq 0, B \zeta \succeq 0, \mathbf{1}^T \zeta = 0\} \quad (1260)$$

From this and (449) we get a vertex-description for a pointed spectral cone that is not full-dimensional:

$$\mathcal{K}_\lambda = \left\{ V_N \begin{pmatrix} \hat{A} \\ \hat{B} \end{pmatrix} V_N^\dagger b \mid b \succeq 0 \right\} \quad (1261)$$

where  $V_N \in \mathbb{R}^{N+1 \times N}$ , and where [*sic*]

$$\hat{B} = e_N^T \in \mathbb{R}^{1 \times N+1} \quad (1262)$$

and

$$\hat{A} = \begin{bmatrix} e_1^T - e_2^T \\ e_2^T - e_3^T \\ \vdots \\ e_{N-1}^T - e_N^T \end{bmatrix} \in \mathbb{R}^{N-1 \times N+1} \quad (1263)$$

hold those rows of  $A$  and  $B$  corresponding to conically independent rows (§2.10) in  $\begin{bmatrix} A \\ B \end{bmatrix} V_N$ .

Conditions (1252) can be equivalently restated in terms of a spectral cone for Euclidean distance matrices:

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} \lambda \begin{pmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{pmatrix} \in \mathcal{K}_M \cap \begin{bmatrix} \mathbb{R}_+^N \\ \mathbb{R}_- \end{bmatrix} \cap \partial \mathcal{H} \\ D \in \mathbb{S}_h^N \end{cases} \quad (1264)$$

Vertex-description of the dual spectral cone is, (319)

$$\begin{aligned} \mathcal{K}_\lambda^* &= \overline{\mathcal{K}_M^* + \begin{bmatrix} \mathbb{R}_+^N \\ \mathbb{R}_- \end{bmatrix}^*} + \partial \mathcal{H}^* \subseteq \mathbb{R}^{N+1} \\ &= \{ [A^T \ B^T \ \mathbf{1} \ -\mathbf{1}] b \mid b \succeq 0 \} = \{ [\hat{A}^T \ \hat{B}^T \ \mathbf{1} \ -\mathbf{1}] a \mid a \succeq 0 \} \end{aligned} \quad (1265)$$

From (1261) and (450) we get a halfspace-description:

$$\mathcal{K}_\lambda^* = \{ y \in \mathbb{R}^{N+1} \mid (V_N^T [\hat{A}^T \ \hat{B}^T])^\dagger V_N^T y \succeq 0 \} \quad (1266)$$

This polyhedral dual spectral cone  $\mathcal{K}_\lambda^*$  is closed, convex, full-dimensional because  $\mathcal{K}_\lambda$  is pointed, but is not pointed because  $\mathcal{K}_\lambda$  is not full-dimensional.

### 5.11.2.3 Unordered eigenspectra

Spectral cones are not unique; eigenspectra ordering can be rendered benign within a cone by presorting a vector of eigenvalues into nonincreasing order.<sup>5.54</sup> Then things simplify: Conditions (1252) now specify eigenvalue membership to the spectral cone

$$\begin{aligned} \lambda \begin{pmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbb{EDM}^N \end{pmatrix} &= \begin{bmatrix} \mathbb{R}_+^N \\ \mathbb{R}_- \end{bmatrix} \cap \partial \mathcal{H} \\ &= \{ \zeta \in \mathbb{R}^{N+1} \mid B \zeta \succeq 0, \ \mathbf{1}^T \zeta = 0 \} \end{aligned} \quad (1267)$$

<sup>5.54</sup>Eigenspectra ordering (represented by a cone having monotone description such as (1256)) becomes benign in (1479), for example, where projection of a given presorted vector on the nonnegative orthant in a subspace is equivalent to its projection on the monotone nonnegative cone in that same subspace; equivalence is a consequence of presorting.

where  $B$  is defined in (1259), and  $\partial\mathcal{H}$  in (1257). From (449) we get a vertex-description for a pointed spectral cone not full-dimensional:

$$\begin{aligned}\lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbb{EDM}^N \end{bmatrix}\right) &= \left\{V_{\mathcal{N}}(\tilde{B}V_{\mathcal{N}})^{\dagger} b \mid b \succeq 0\right\} \\ &= \left\{\begin{bmatrix} I \\ -\mathbf{1}^T \end{bmatrix} b \mid b \succeq 0\right\}\end{aligned}\quad (1268)$$

where  $V_{\mathcal{N}} \in \mathbb{R}^{N+1 \times N}$  and

$$\tilde{B} \triangleq \begin{bmatrix} e_1^T \\ e_2^T \\ \vdots \\ e_N^T \end{bmatrix} \in \mathbb{R}^{N \times N+1} \quad (1269)$$

holds only those rows of  $B$  corresponding to conically independent rows in  $BV_{\mathcal{N}}$ .

For presorted eigenvalues, (1252) can be equivalently restated

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} \lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix}\right) \in \left[\begin{array}{c} \mathbb{R}_+^N \\ \mathbb{R}_- \end{array}\right] \cap \partial\mathcal{H} \\ D \in \mathbb{S}_h^N \end{cases} \quad (1270)$$

Vertex-description of the dual spectral cone is, (319)

$$\begin{aligned}\lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbb{EDM}^N \end{bmatrix}\right)^* &= \left[\begin{array}{c} \mathbb{R}_+^N \\ \mathbb{R}_- \end{array}\right] + \partial\mathcal{H}^* \subseteq \mathbb{R}^{N+1} \\ &= \{[B^T \ \mathbf{1} \ -\mathbf{1}] b \mid b \succeq 0\} = \{\tilde{B}^T \ \mathbf{1} \ -\mathbf{1} a \mid a \succeq 0\}\end{aligned}\quad (1271)$$

From (450) we get a halfspace-description:

$$\begin{aligned}\lambda\left(\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbb{EDM}^N \end{bmatrix}\right)^* &= \{y \in \mathbb{R}^{N+1} \mid (V_{\mathcal{N}}^T \tilde{B}^T)^{\dagger} V_{\mathcal{N}}^T y \succeq 0\} \\ &= \{y \in \mathbb{R}^{N+1} \mid [I \ -\mathbf{1}] y \succeq 0\}\end{aligned}\quad (1272)$$

This polyhedral dual spectral cone is closed, convex, full-dimensional but not pointed.<sup>5.55</sup>

#### 5.11.2.4 Dual cone *versus* dual spectral cone

An open question regards the relationship of convex cones and their duals to the corresponding spectral cones and their duals. A positive semidefinite cone, for example, is selfdual. Both the nonnegative orthant and the monotone nonnegative cone are spectral cones for it. When we consider the nonnegative orthant, then that spectral cone for the selfdual positive semidefinite cone is also selfdual.

## 5.12 List reconstruction

The term *metric multidimensional scaling*<sup>5.56</sup> [287] [118] [394] [116] [291] [97] refers to any reconstruction of a list  $X \in \mathbb{R}^{n \times N}$  in Euclidean space from interpoint distance information, possibly incomplete (§6.7), ordinal (§5.13.2), or specified perhaps only by

<sup>5.55</sup>Notice that any nonincreasingly ordered eigenspectrum belongs to this dual spectral cone.

<sup>5.56</sup>Scaling [390] means making a scale, i.e., a numerical representation of qualitative data. If the scale is multidimensional, it's multidimensional scaling.

—Jan de Leeuw

In one dimension,  $N$  coordinates in  $X$  define the scale; e.g., §7.2.2.7.1.

bounding-constraints (§5.4.2.2.12) [392]. Techniques for reconstruction are essentially methods for optimally embedding an unknown list of points, corresponding to given Euclidean distance data, in an affine subset of desired or minimum dimension. The oldest known precursor is called *principal component analysis* [189] which analyzes the correlation matrix (§5.9.1.0.1); [56, §22] a.k.a, *Karhunen-Loéve transform* in digital signal processing literature.

A goal of multidimensional scaling is to find a low-dimensional representation of list  $X$  so that distances between its elements best preserve a given set of pairwise dissimilarities. *Dissimilarity* is some measure or perception of unlikeness. *Similarity* between vectors (in Euclidean space) is measured by inner product, [328, §2] whereas dissimilarity is measured by distance-square.<sup>5.57</sup> When dissimilarity data comprises measurable distances, then reconstruction is termed *metric* multidimensional scaling.

Isometric reconstruction (§5.5.3) of point list  $X$  is best performed by eigenvalue decomposition of a Gram matrix; for then, numerical errors of factorization are easily spotted in the eigenvalues: Now we consider how rotation/reflection and translation invariance factor into a reconstruction.

### 5.12.1 $x_1$ at the origin. $V_N$

At the stage of reconstruction, we have  $D \in \mathbb{EDM}^N$  and wish to find a generating list (§2.3.2) for polyhedron  $\mathcal{P} - \alpha$  by factoring Gram matrix  $-V_N^T DV_N$  (1050) as in (1235). One way to factor  $-V_N^T DV_N$  is via diagonalization of symmetric matrices; [368, §5.6] [228] (§A.5.1, §A.3)

$$-V_N^T DV_N \triangleq Q \Lambda Q^T \quad (1273)$$

$$Q \Lambda Q^T \succeq 0 \Leftrightarrow \Lambda \succeq 0 \quad (1274)$$

where  $Q \in \mathbb{R}^{N-1 \times N-1}$  is an orthogonal matrix containing eigenvectors while  $\Lambda \in \mathbb{S}^{N-1}$  is a diagonal matrix containing corresponding nonnegative eigenvalues ordered by nonincreasing value. From the diagonalization, identify the list using (1180);

$$-V_N^T DV_N = 2V_N^T X^T X V_N \triangleq Q \sqrt{\Lambda} Q_p^T Q \sqrt{\Lambda} Q^T \quad (1275)$$

where  $\sqrt{\Lambda} Q_p^T Q \sqrt{\Lambda} \triangleq \Lambda = \sqrt{\Lambda} \sqrt{\Lambda}$  and where  $Q_p \in \mathbb{R}^{n \times N-1}$  is unknown as is its dimension  $n$ . Rotation/reflection is accounted for by  $Q_p$  yet only its first  $r$  columns are necessarily orthonormal.<sup>5.58</sup> Assuming membership to the unit simplex  $y \in \mathcal{S}$  (1232), then point  $p = X\sqrt{2}V_N y = Q_p \sqrt{\Lambda} Q^T y$  in  $\mathbb{R}^n$  belongs to the translated polyhedron

$$\mathcal{P} - x_1 \quad (1276)$$

whose generating list constitutes the columns of (1174)

$$\begin{aligned} [\mathbf{0} \quad X\sqrt{2}V_N] &= [\mathbf{0} \quad Q_p \sqrt{\Lambda} Q^T] \in \mathbb{R}^{n \times N} \\ &= [\mathbf{0} \quad x_2 - x_1 \quad x_3 - x_1 \quad \cdots \quad x_N - x_1] \end{aligned} \quad (1277)$$

The scaled auxiliary matrix  $V_N$  represents that translation. A simple choice for  $Q_p$  has  $n$  set to  $N-1$ ; *id est*,  $Q_p = I$ . Ideally, each member of the generating list has at most  $r$  nonzero entries;  $r$  being, affine dimension

$$\text{rank } V_N^T DV_N = \text{rank } Q_p \sqrt{\Lambda} Q^T = \text{rank } \Lambda = r \quad (1278)$$

<sup>5.57</sup>This sense reversal is analogous to autocorrelation *versus* total power of lagged differences in digital signal processing. [105, p.9]

<sup>5.58</sup>Recall  $r$  signifies affine dimension.  $Q_p$  is not necessarily an orthogonal matrix.  $Q_p$  is constrained such that only its first  $r$  columns are necessarily orthonormal because there are only  $r$  nonzero eigenvalues in  $\Lambda$  when  $-V_N^T DV_N$  has rank  $r$  (§5.7.1.1). Remaining columns of  $Q_p$  are arbitrary.

Each member then has at least  $N-1-r$  zeros in its higher-dimensional coordinates because  $r \leq N-1$ . (1186) To truncate those zeros, choose  $n$  equal to affine dimension which is the smallest  $n$  possible because  $XV_N$  has rank  $r \leq n$  (1182).<sup>5.59</sup> In that case, the simplest choice for  $Q_p$  is  $[I \ \mathbf{0}]$  having dimension  $r \times N-1$ .

We may wish to verify the list (1277) found from the diagonalization of  $-V_N^T DV_N$ . Because of rotation/reflection and translation invariance (§5.5), EDM  $D$  can be uniquely made from that list by calculating: (1033)

$$\mathbf{D}(X) = \mathbf{D}(X[\mathbf{0} \ \sqrt{2}V_N]) = \mathbf{D}(Q_p[\mathbf{0} \ \sqrt{\Lambda} Q^T]) = \mathbf{D}([\mathbf{0} \ \sqrt{\Lambda} Q^T]) \quad (1279)$$

This suggests a way to find EDM  $D$  given  $-V_N^T DV_N$  (*confer* (1158))

$$D = \begin{bmatrix} 0 \\ \delta(-V_N^T DV_N) \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & \delta(-V_N^T DV_N)^T \end{bmatrix} - 2 \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & -V_N^T DV_N \end{bmatrix} \quad (1154)$$

### 5.12.2 0 geometric center. $V$

Alternatively we may perform reconstruction using auxiliary matrix  $V$  (§B.4.1) and Gram matrix  $-VDV\frac{1}{2}$  (1054) instead; to find a generating list for polyhedron

$$\mathcal{P} - \alpha_c \quad (1280)$$

whose geometric center  $\alpha_c$  has been translated to the origin. Redimensioning diagonalization factors  $Q, \Lambda \in \mathbb{R}^{N \times N}$  and unknown  $Q_p \in \mathbb{R}^{n \times N}$ , (1181)

$$-VDV = 2VX^T XV \triangleq Q\Lambda Q^T \triangleq Q\sqrt{\Lambda}Q_p^T Q_p \sqrt{\Lambda} Q^T \quad (1281)$$

where the geometrically centered generating list constitutes (*confer* (1277))

$$\begin{aligned} XV &= \frac{1}{\sqrt{2}} Q_p \sqrt{\Lambda} Q^T \in \mathbb{R}^{n \times N} \\ &= [x_1 - \frac{1}{N} X \mathbf{1} \quad x_2 - \frac{1}{N} X \mathbf{1} \quad x_3 - \frac{1}{N} X \mathbf{1} \quad \cdots \quad x_N - \frac{1}{N} X \mathbf{1}] \end{aligned} \quad (1282)$$

where  $\alpha_c = \frac{1}{N} X \mathbf{1}$ . (§5.5.1.0.1) Recall,  $Q_p$  accounts for list rotation/reflection. The simplest choice for  $Q_p$  is  $[I \ \mathbf{0}] \in \mathbb{R}^{r \times N}$  with affine dimension  $r$ .

Now EDM  $D$  can be uniquely made from the list found: (1033)

$$\mathbf{D}(X) = \mathbf{D}(XV) = \mathbf{D}(\frac{1}{\sqrt{2}} Q_p \sqrt{\Lambda} Q^T) = \mathbf{D}(\sqrt{\Lambda} Q^T) \frac{1}{2} \quad (1283)$$

This EDM is, of course, identical to (1279). Similarly to (1154), from  $-VDV$  we can find EDM  $D$  (*confer* (1145))

$$D = \delta(-VDV\frac{1}{2})\mathbf{1}^T + \mathbf{1}\delta(-VDV\frac{1}{2})^T - 2(-VDV\frac{1}{2}) \quad (1144)$$

---

<sup>5.59</sup>If we write  $Q^T = \begin{bmatrix} q_1^T \\ \vdots \\ q_{N-1}^T \end{bmatrix}$  as rowwise eigenvectors,  $\Lambda = \begin{bmatrix} \lambda_1 & & & \mathbf{0} \\ & \ddots & & \\ & & \lambda_r & \\ \mathbf{0} & & & \ddots & & \mathbf{0} \end{bmatrix}$  in terms of eigenvalues,

and  $Q_p = [q_{p1} \cdots q_{pN-1}]$  as column vectors, then  $Q_p \sqrt{\Lambda} Q^T = \sum_{i=1}^r \sqrt{\lambda_i} q_{pi} q_i^T$  is a sum of  $r$  linearly independent rank-one matrices (§B.1.1). Hence the summation has rank  $r$ .

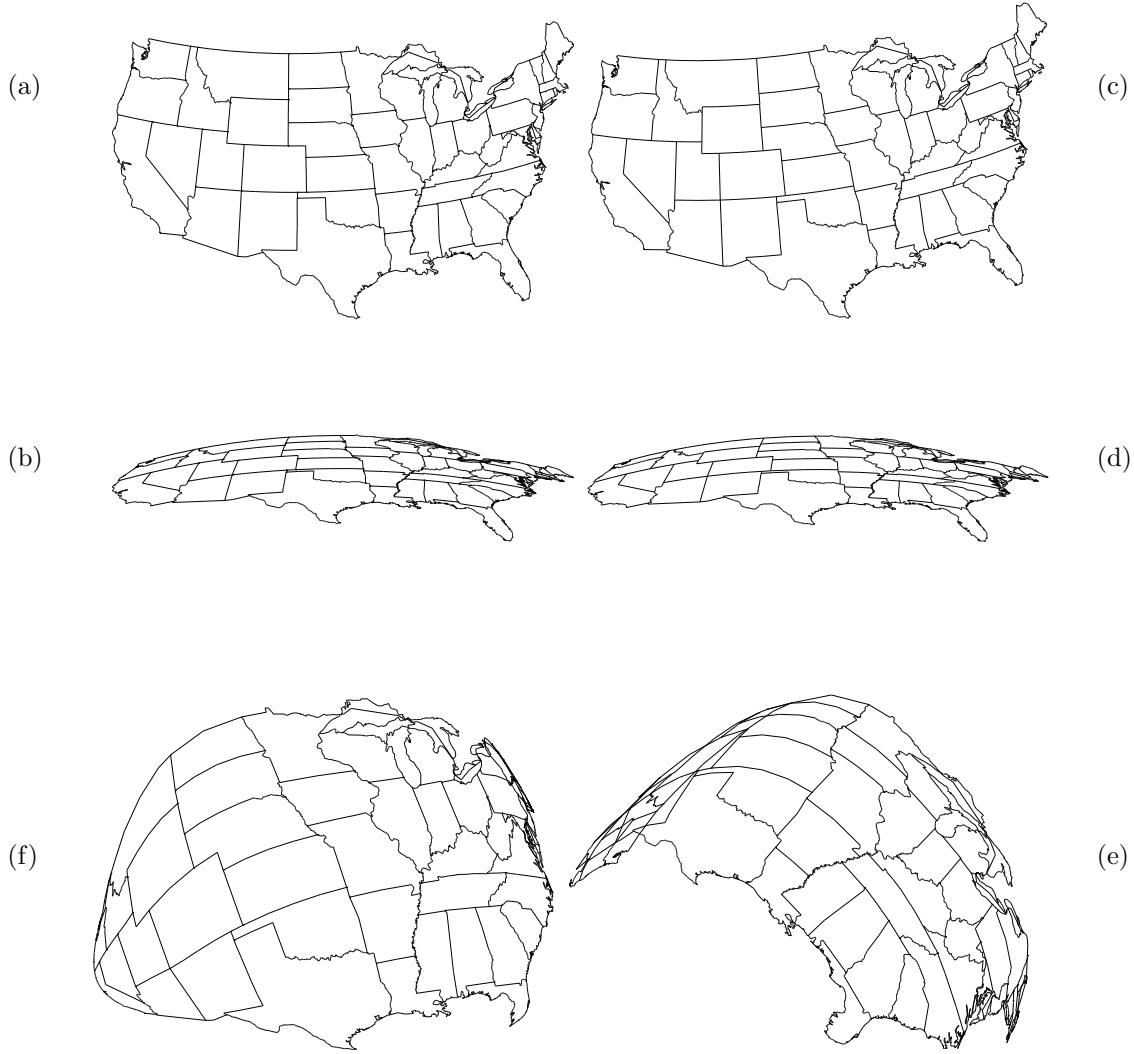


Figure 159: (*confer Figure 8*) Nonconvex map of United States of America showing some state boundaries and the Great Lakes. All plots made by connecting 5020 points. Any difference in scale in (a) through (d) is artifact of plotting routine.

- (a) Shows original map made from decimated (latitude, longitude) data.
- (b) Original map data rotated (freehand) to highlight curvature of Earth.
- (c) Map isometrically reconstructed from an EDM (from distance only).
- (d) Same reconstructed map illustrating curvature.
- (e)(f) Two views of one isotonic reconstruction (from comparative distance); problem (1293) with no sort constraint  $\Pi \underline{d}$  (and no hidden line removal).

## 5.13 Reconstruction examples

### 5.13.1 Isometric reconstruction

#### 5.13.1.0.1 Example. Cartography.

The most fundamental application of EDMs is to reconstruct relative point position given only interpoint distance information. Drawing a map of the United States is a good illustration of isometric reconstruction (§5.4.2.2.10) from complete distance data. We obtained latitude and longitude information for the coast, border, states, and Great Lakes from the [usalo atlas data file](#) within MATLAB Mapping Toolbox; conversion to Cartesian coordinates  $(x, y, z)$  via:

$$\begin{aligned}\phi &\triangleq \pi/2 - \text{latitude} \\ \theta &\triangleq \text{longitude} \\ x &= \sin(\phi) \cos(\theta) \\ y &= \sin(\phi) \sin(\theta) \\ z &= \cos(\phi)\end{aligned}\tag{1284}$$

We used 64% of the available map data to calculate EDM  $D$  from  $N = 5020$  points. The original (decimated) data and its isometric reconstruction via (1275) are shown in Figure 159a-d. [423, MATLAB code] The eigenvalues computed for (1273) are

$$\lambda(-V_N^T D V_N) = [199.8 \ 152.3 \ 2.465 \ 0 \ 0 \ 0 \ \dots]^T\tag{1285}$$

The 0 eigenvalues have absolute numerical error on the order of 2E-13; meaning, the EDM data indicates three dimensions ( $r = 3$ ) are required for reconstruction to nearly machine precision.  $\square$

### 5.13.2 Isotonic reconstruction

Sometimes only comparative information about distance is known (Earth is closer to the Moon than it is to the Sun). Suppose, for example, EDM  $D$  for three points is unknown:

$$D = [d_{ij}] = \begin{bmatrix} 0 & d_{12} & d_{13} \\ d_{12} & 0 & d_{23} \\ d_{13} & d_{23} & 0 \end{bmatrix} \in \mathbb{S}_h^3\tag{1022}$$

but comparative distance data is available:

$$d_{13} \geq d_{23} \geq d_{12}\tag{1286}$$

With vectorization  $\underline{d} = [d_{12} \ d_{13} \ d_{23}]^T \in \mathbb{R}^3$ , we express the comparative data as the nonincreasing sorting

$$\Pi \underline{d} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} d_{12} \\ d_{13} \\ d_{23} \end{bmatrix} = \begin{bmatrix} d_{13} \\ d_{23} \\ d_{12} \end{bmatrix} \in \mathcal{K}_{\mathcal{M}+}\tag{1287}$$

where  $\Pi$  is a given permutation matrix expressing known sorting action on the entries of unknown EDM  $D$ , and  $\mathcal{K}_{\mathcal{M}+}$  is the monotone nonnegative cone (§2.13.10.4.2)

$$\mathcal{K}_{\mathcal{M}+} = \{z \mid z_1 \geq z_2 \geq \dots \geq z_{N(N-1)/2} \geq 0\} \subseteq \mathbb{R}_+^{N(N-1)/2}\tag{434}$$

where  $N(N-1)/2 = 3$  for the present example. From sorted vectorization (1287) we create the *sort-index matrix*

$$O = \begin{bmatrix} 0 & 1^2 & 3^2 \\ 1^2 & 0 & 2^2 \\ 3^2 & 2^2 & 0 \end{bmatrix} \in \mathbb{S}_h^3 \cap \mathbb{R}_+^{3 \times 3}\tag{1288}$$

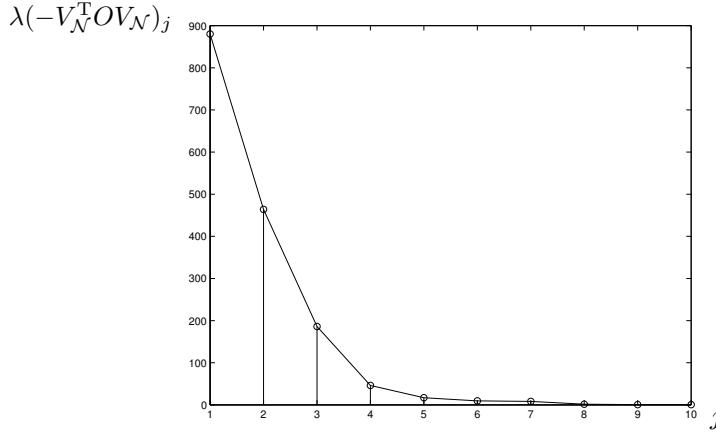


Figure 160: Largest ten eigenvalues, of  $-V_N^T O V_N$  for USA map, sorted by decreasing value.

generally defined

$$O_{ij} \triangleq k^2 \mid d_{ij} = (\Xi \Pi \underline{d})_k , \quad j \neq i \quad (1289)$$

where  $\Xi$  is a permutation matrix (1900) completely reversing order of vector entries.

Replacing EDM data with indices-square of a nonincreasing sorting like this is, of course, a heuristic we invented and may be regarded as a nonlinear introduction of much noise into the Euclidean distance matrix. For large data sets, this heuristic makes an otherwise intense problem computationally tractable; we see an example in relaxed problem (1294).

Any process of reconstruction that leaves comparative distance information intact is called *ordinal multidimensional scaling* or *isotonic reconstruction*. Beyond rotation, reflection, and translation error, (§5.5) list reconstruction by isotonic reconstruction is subject to error in absolute scale (*dilation*) and distance ratio. Yet Borg & Groenen argue: [56, §2.2] reconstruction from complete comparative distance information for a large number of points is as highly constrained as reconstruction from an EDM; the larger the number, the smaller the optimal solution set; whereas,

$$\text{isotonic solution set} \supseteq \text{isometric solution set} \quad (1290)$$

### 5.13.2.1 Isotonic cartography

To test Borg & Groenen's conjecture, suppose we make a complete sort-index matrix  $O \in \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N}$  for the map of USA and then substitute  $O$  in place of EDM  $D$  in the reconstruction process of §5.12. Whereas EDM  $D$  returned only three significant eigenvalues (1285), the sort-index matrix  $O$  is generally not an EDM (certainly not an EDM with corresponding affine dimension 3) so returns many more. The eigenvalues, calculated with absolute numerical error approximately 5E-7, are plotted in Figure 160:

$$\lambda(-V_N^T O V_N) = [880.1 \ 463.9 \ 186.1 \ 46.20 \ 17.12 \ 9.625 \ 8.257 \ 1.701 \ 0.7128 \ 0.6460 \dots]^T \quad (1291)$$

The extra eigenvalues indicate that affine dimension corresponding to an EDM near  $O$  is likely to exceed 3. To realize the map, we must simultaneously reduce that dimensionality and find an EDM  $D$  closest to  $O$  in some sense<sup>5.60</sup> while maintaining

---

<sup>5.60</sup> a problem explored more in §7.

the known comparative distance relationship. For example: given permutation matrix  $\Pi$  expressing the known sorting action like (1287) on entries

$$\underline{d} \triangleq \frac{1}{\sqrt{2}} \text{dvec } D = \begin{bmatrix} d_{12} \\ d_{13} \\ d_{23} \\ d_{14} \\ d_{24} \\ d_{34} \\ \vdots \\ d_{N-1,N} \end{bmatrix} \in \mathbb{R}^{N(N-1)/2} \quad (1292)$$

of unknown  $D \in \mathbb{S}_h^N$ , we can make sort-index matrix  $O$  input to the optimization problem

$$\begin{array}{ll} \underset{D}{\text{minimize}} & \| -V_N^T(D - O)V_N \|_F \\ \text{subject to} & \text{rank } V_N^T D V_N \leq 3 \\ & \Pi \underline{d} \in \mathcal{K}_{\mathcal{M}+} \\ & D \in \mathbb{EDM}^N \end{array} \quad (1293)$$

that finds the EDM  $D$  (corresponding to affine dimension not exceeding 3 in isomorphic  $\text{dvec } \mathbb{EDM}^N \cap \Pi^T \mathcal{K}_{\mathcal{M}+}$ ) closest to  $O$  in the sense of Schoenberg (1052).

Analytical solution to this problem, ignoring the sort constraint  $\Pi \underline{d} \in \mathcal{K}_{\mathcal{M}+}$ , is known [394]: we get the convex optimization [*sic*] (§7.1)

$$\begin{array}{ll} \underset{D}{\text{minimize}} & \| -V_N^T(D - O)V_N \|_F \\ \text{subject to} & \text{rank } V_N^T D V_N \leq 3 \\ & D \in \mathbb{EDM}^N \end{array} \quad (1294)$$

Only the three largest nonnegative eigenvalues in (1291) need be retained to make list (1277); the rest are discarded. The reconstruction from EDM  $D$  found in this manner is plotted in Figure 159e-f. (In the MATLAB code on *Wukimization* [416], matrix  $O$  is normalized by  $(N(N-1)/2)^2$ .) From these plots it becomes obvious that inclusion of the sort constraint is necessary for isotonic reconstruction.

That sort constraint demands: any optimal solution  $D^*$  must possess the known comparative distance relationship that produces the original ordinal distance data  $O$  (1289). Ignoring the sort constraint, apparently, violates it. Yet even more remarkable is how much the map, reconstructed using only ordinal data, still resembles the original map of USA after suffering the many violations produced by solving relaxed problem (1294). This suggests the simple reconstruction techniques of §5.12 are robust to a significant amount of noise.

### 5.13.2.2 Isotonic solution with sort constraint

Because problems involving rank are generally difficult, we will partition (1293) into two problems we know how to solve and then alternate their solution until convergence:

$$\begin{array}{ll} \underset{D}{\text{minimize}} & \| -V_N^T(D - O)V_N \|_F \\ \text{subject to} & \text{rank } V_N^T D V_N \leq 3 \\ & D \in \mathbb{EDM}^N \end{array} \quad (\text{a}) \quad (1294)$$

$$\begin{array}{ll} \underset{\sigma}{\text{minimize}} & \|\sigma - \Pi \underline{d}\| \\ \text{subject to} & \sigma \in \mathcal{K}_{\mathcal{M}+} \end{array} \quad (\text{b}) \quad (1295)$$

where sort-index matrix  $O$  (a given constant in (a)) becomes an implicit vector variable  $\underline{o}_i$  solving the  $i^{\text{th}}$  instance of (1295b)

$$\frac{1}{\sqrt{2}} \text{dvec } O_i = \underline{o}_i \triangleq \Pi^T \sigma^* \in \mathbb{R}^{N(N-1)/2}, \quad i \in \{1, 2, 3, \dots\} \quad (1296)$$

As mentioned in discussion of relaxed problem (1294), a closed-form solution to problem (1295a) exists. Only the first iteration of (1295a) sees the original sort-index matrix  $O$  whose entries are nonnegative whole numbers; *id est*,  $O_0 = O \in \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N}$  (1289). Subsequent iterations  $i$  take the previous solution of (1295b) as input

$$O_i = \text{dvec}^{-1}(\sqrt{2} \underline{o}_i) \in \mathbb{S}^N \quad (1297)$$

real successors, estimating distance-square not order, to the sort-index matrix  $O$ .

New convex problem (1295b) finds the unique minimum-distance projection of  $\Pi \underline{d}$  on the monotone nonnegative cone  $\mathcal{K}_{\mathcal{M}+}$ . By defining

$$Y^{\dagger T} = [e_1 - e_2 \quad e_2 - e_3 \quad e_3 - e_4 \quad \cdots \quad e_m] \in \mathbb{R}^{m \times m} \quad (435)$$

where  $m \triangleq N(N-1)/2$ , we may rewrite (1295b) as an equivalent quadratic program; a convex problem in terms of the halfspace-description of  $\mathcal{K}_{\mathcal{M}+}$ :

$$\begin{aligned} & \underset{\sigma}{\text{minimize}} \quad (\sigma - \Pi \underline{d})^T (\sigma - \Pi \underline{d}) \\ & \text{subject to} \quad Y^{\dagger} \sigma \succeq 0 \end{aligned} \quad (1298)$$

This quadratic program can be converted to a semidefinite program via Schur-form (§3.5.3); we get the equivalent problem

$$\begin{aligned} & \underset{t \in \mathbb{R}, \sigma}{\text{minimize}} \quad t \\ & \text{subject to} \quad \begin{bmatrix} tI & \sigma - \Pi \underline{d} \\ (\sigma - \Pi \underline{d})^T & 1 \end{bmatrix} \succeq 0 \\ & \quad Y^{\dagger} \sigma \succeq 0 \end{aligned} \quad (1299)$$

### 5.13.2.3 Convergence

In §E.10 we discuss convergence of alternating projection on intersecting convex sets in a Euclidean vector space; convergence to a point in their intersection. Here the situation is different for two reasons:

Firstly, sets of positive semidefinite matrices having an upper bound on rank are generally not convex. Yet in §7.1.4.0.1 we prove that (1295a) is equivalent to a projection of nonincreasingly ordered eigenvalues on a subset of the nonnegative orthant:

$$\begin{aligned} & \underset{D}{\text{minimize}} \quad \| -V_{\mathcal{N}}^T (D - O) V_{\mathcal{N}} \|_F \quad & & \underset{\Upsilon}{\text{minimize}} \quad \| \Upsilon - \Lambda \|_F \\ & \text{subject to} \quad \text{rank } V_{\mathcal{N}}^T D V_{\mathcal{N}} \leq 3 \quad \equiv \quad \text{subject to} \quad \delta(\Upsilon) \in \begin{bmatrix} \mathbb{R}_+^3 \\ \mathbf{0} \end{bmatrix} \end{aligned} \quad (1300)$$

where  $-V_{\mathcal{N}}^T D V_{\mathcal{N}} \triangleq U \Upsilon U^T \in \mathbb{S}^{N-1}$  and  $-V_{\mathcal{N}}^T O V_{\mathcal{N}} \triangleq Q \Lambda Q^T \in \mathbb{S}^{N-1}$  are ordered diagonalizations (§A.5). It so happens: optimal orthogonal  $U^*$  always equals  $Q$  given. Linear operator  $T(A) = U^{*\top} A U^*$ , acting on square matrix  $A$ , is an isometry because Frobenius' norm is orthogonally invariant (49). This isometric isomorphism  $T$  thus maps a nonconvex problem to a convex one that preserves distance.

Secondly, the second half (1295b) of the *alternation* takes place in a different vector space;  $\mathbb{S}_h^N$  (*versus*  $\mathbb{S}^{N-1}$ ). From §5.6 we know these two vector spaces are related by an isomorphism,  $\mathbb{S}^{N-1} = \mathbf{V}_{\mathcal{N}}(\mathbb{S}_h^N)$  (1163), but not by an isometry.

We have, therefore, no guarantee from theory of alternating projection that alternation (1295) converges to a point, in the set of all EDMs corresponding to affine dimension not in excess of 3, belonging to  $\text{dvec } \mathbb{EDM}^N \cap \Pi^T \mathcal{K}_{\mathcal{M}+}$ .

### 5.13.2.4 Interlude

Map reconstruction from comparative distance data, isotonic reconstruction, would also prove invaluable to stellar cartography where absolute interstellar distance is difficult to acquire. But we have not yet implemented the second half (1298) of alternation (1295) for USA map data because memory-demands exceed capability of our computer.

#### 5.13.2.4.1 Exercise. *Convergence of isotonic solution by alternation.*

Empirically demonstrate convergence, discussed in §5.13.2.3, on a smaller data set. ▼

It would be remiss not to mention another method of solution to this isotonic reconstruction problem: Once again we assume only comparative distance data like (1286) is available. Given known set of indices  $\mathcal{I}$

$$\begin{aligned} & \underset{D}{\text{minimize}} \quad \text{rank } VDV \\ & \text{subject to } d_{ij} \leq d_{kl} \leq d_{mn} \quad \forall (i, j, k, l, m, n) \in \mathcal{I} \\ & \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1301)$$

this problem minimizes affine dimension while finding an EDM whose entries satisfy known comparative relationships. Suitable rank heuristics are discussed in §4.5.1 and §7.2.2 that will transform this to a convex optimization problem.

Using contemporary computers, even with a rank heuristic in place of the objective function, this problem formulation is more difficult to compute than the relaxed counterpart problem (1294). That is because there exist efficient algorithms to compute a selected few eigenvalues and eigenvectors from a very large matrix. Regardless, it is important to recognize: the optimal solution set for this problem (1301) is practically always different from the optimal solution set for its counterpart, problem (1293).

## 5.14 Fifth property of Euclidean metric

We continue now with the question raised in §5.3 regarding necessity for at least one requirement more than the four properties of the Euclidean metric (§5.2) to certify realizability of a bounded convex polyhedron or to reconstruct a generating list for it from incomplete distance information. There we saw that the four Euclidean metric properties are necessary for  $D \in \mathbb{EDM}^N$  in the case  $N=3$ , but become insufficient when cardinality  $N$  exceeds 3 (regardless of affine dimension).

### 5.14.1 Recapitulate

In the particular case  $N=3$ ,  $-V_N^T DV_N \succeq 0$  (1204) and  $D \in \mathbb{S}_h^3$  are necessary and sufficient conditions for  $D$  to be an EDM. By (1206), triangle inequality is then the only Euclidean condition bounding the necessarily nonnegative  $d_{ij}$ ; and those bounds are tight. That means the first four properties of the Euclidean metric are necessary and sufficient conditions for  $D$  to be an EDM in the case  $N=3$ ; for  $i, j \in \{1, 2, 3\}$

$$\begin{aligned} \sqrt{d_{ij}} &\geq 0, \quad i \neq j \\ \sqrt{d_{ij}} &= 0, \quad i = j \\ \sqrt{d_{ij}} &= \sqrt{d_{ji}} \\ \sqrt{d_{ij}} &\leq \sqrt{d_{ik}} + \sqrt{d_{kj}}, \quad i \neq j \neq k \end{aligned} \Leftrightarrow \begin{aligned} -V_N^T DV_N &\succeq 0 \\ D &\in \mathbb{S}_h^3 \end{aligned} \Leftrightarrow D \in \mathbb{EDM}^3 \quad (1302)$$

Yet those four properties become insufficient when  $N > 3$ .

### 5.14.2 Derivation of the Fifth

Correspondence between the triangle inequality and the EDM was developed in §5.8.2 where a triangle inequality (1206a) was revealed within the leading principal  $2 \times 2$  submatrix of  $-V_N^T D V_N$  when positive semidefinite. Our choice of the leading principal submatrix was arbitrary; actually, a unique triangle inequality like (1101) corresponds to any one of the  $(N-1)!/(2!(N-1-2)!)$  principal  $2 \times 2$  submatrices.<sup>5.61</sup> Assuming  $D \in \mathbb{S}_h^4$  and  $-V_N^T D V_N \in \mathbb{S}^3$ , then by the *positive (semi)definite principal submatrices theorem* (§A.3.1.0.4) it is sufficient to prove: all  $d_{ij}$  are nonnegative, all triangle inequalities are satisfied, and  $\det(-V_N^T D V_N)$  is nonnegative. When  $N=4$ , in other words, that nonnegative determinant becomes the fifth and last Euclidean metric requirement for  $D \in \mathbb{EDM}^N$ . We now endeavor to ascribe geometric meaning to it.

#### 5.14.2.1 Nonnegative determinant

By (1107) when  $D \in \mathbb{EDM}^4$ ,  $-V_N^T D V_N$  is equal to inner product (1102),

$$\Theta^T \Theta = \begin{bmatrix} d_{12} & \sqrt{d_{12}d_{13}} \cos \theta_{213} & \sqrt{d_{12}d_{14}} \cos \theta_{214} \\ \sqrt{d_{12}d_{13}} \cos \theta_{213} & d_{13} & \sqrt{d_{13}d_{14}} \cos \theta_{314} \\ \sqrt{d_{12}d_{14}} \cos \theta_{214} & \sqrt{d_{13}d_{14}} \cos \theta_{314} & d_{14} \end{bmatrix} \quad (1303)$$

Because Euclidean space is an inner-product space, the more concise inner-product form of the determinant is admitted;

$$\det(\Theta^T \Theta) = -d_{12}d_{13}d_{14}(\cos(\theta_{213})^2 + \cos(\theta_{214})^2 + \cos(\theta_{314})^2 - 2 \cos \theta_{213} \cos \theta_{214} \cos \theta_{314} - 1) \quad (1304)$$

The determinant is nonnegative if and only if

$$\begin{aligned} \cos \theta_{214} \cos \theta_{314} - \sqrt{\sin(\theta_{214})^2 \sin(\theta_{314})^2} &\leq \cos \theta_{213} \leq \cos \theta_{214} \cos \theta_{314} + \sqrt{\sin(\theta_{214})^2 \sin(\theta_{314})^2} \\ &\Leftrightarrow \\ \cos \theta_{213} \cos \theta_{314} - \sqrt{\sin(\theta_{213})^2 \sin(\theta_{314})^2} &\leq \cos \theta_{214} \leq \cos \theta_{213} \cos \theta_{314} + \sqrt{\sin(\theta_{213})^2 \sin(\theta_{314})^2} \\ &\Leftrightarrow \\ \cos \theta_{213} \cos \theta_{214} - \sqrt{\sin(\theta_{213})^2 \sin(\theta_{214})^2} &\leq \cos \theta_{314} \leq \cos \theta_{213} \cos \theta_{214} + \sqrt{\sin(\theta_{213})^2 \sin(\theta_{214})^2} \end{aligned} \quad (1305)$$

which simplifies, for  $0 \leq \theta_{i1\ell}, \theta_{\ell 1j}, \theta_{i1j} \leq \pi$  and all  $i \neq j \neq \ell \in \{2, 3, 4\}$ , to

$$\cos(\theta_{i1\ell} + \theta_{\ell 1j}) \leq \cos \theta_{i1j} \leq \cos(\theta_{i1\ell} - \theta_{\ell 1j}) \quad (1306)$$

Analogously to triangle inequality (1218), the determinant is 0 upon equality on either side of (1306) which is tight. Inequality (1306) can be equivalently written linearly as a triangle inequality between relative angles [455, §1.4];

$$\begin{aligned} |\theta_{i1\ell} - \theta_{\ell 1j}| &\leq \theta_{i1j} \leq \theta_{i1\ell} + \theta_{\ell 1j} \\ \theta_{i1\ell} + \theta_{\ell 1j} + \theta_{i1j} &\leq 2\pi \\ 0 &\leq \theta_{i1\ell}, \theta_{\ell 1j}, \theta_{i1j} \leq \pi \end{aligned} \quad (1307)$$

Generalizing this:

---

<sup>5.61</sup>There are fewer principal  $2 \times 2$  submatrices in  $-V_N^T D V_N$  than there are triangles made by four or more points because there are  $N!/(3!(N-3)!)$  triangles made by point triples. The triangles corresponding to those submatrices all have vertex  $x_1$ . (confer §5.8.2.1)

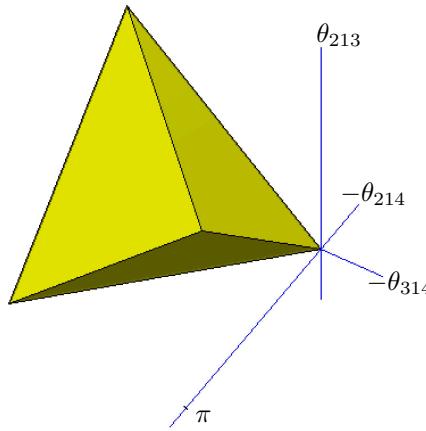


Figure 161: The *relative-angle inequality tetrahedron* (1308) bounding  $\mathbb{EDM}^4$  is regular; drawn in entirety. Each angle  $\theta$  (1099) must belong to this solid to be realizable.

#### 5.14.2.1.1 Fifth property of Euclidean metric - restatement.

*Relative-angle inequality.*

[53] [54, p.17, p.107] [264, §3.1]

(confer §5.3.1.0.1) Augmenting the four fundamental Euclidean metric properties in  $\mathbb{R}^n$ , for all  $i, j, \ell \neq k \in \{1 \dots N\}$ ,  $i < j < \ell$ , and for  $N \geq 4$  distinct points  $\{x_k\}$ , the inequalities

$$\begin{aligned} |\theta_{ik\ell} - \theta_{\ellkj}| &\leq \theta_{ikj} \leq \theta_{ik\ell} + \theta_{\ellkj} & (a) \\ \theta_{ik\ell} + \theta_{\ellkj} + \theta_{ikj} &\leq 2\pi & (b) \\ 0 &\leq \theta_{ik\ell}, \theta_{\ellkj}, \theta_{ikj} \leq \pi & (c) \end{aligned} \quad (1308)$$

must be satisfied at each point  $x_k$  regardless of affine dimension, where  $\theta_{ikj} = \theta_{jki}$  is the angle between vectors at vertex  $x_k$  as defined in (1099) and illustrated in Figure 142.

◊

Because point labelling is arbitrary, this fifth Euclidean metric requirement must apply to each of the  $N$  points as though each were in turn labelled  $x_1$ ; hence the new index  $k$  in (1308). Just as the triangle inequality is the ultimate test for realizability of only three points, the relative-angle inequality is the ultimate test for only four. For four distinct points, the triangle inequality remains a necessary although penultimate test; (§5.4.3)

$$\text{Four Euclidean metric properties (§5.2). } \Leftrightarrow -V_N^T D V_N \succeq 0 \Leftrightarrow D = \mathbf{D}(\Theta) \in \mathbb{EDM}^4 \quad (1309)$$

Angle  $\theta$  inequality (1027) or (1308).  $D \in \mathbb{S}_h^4$

The relative-angle inequality, for this case, is illustrated in Figure 161.

#### 5.14.2.2 Beyond the fifth metric property

When cardinality  $N$  exceeds 4, the first four properties of the Euclidean metric and the relative-angle inequality together become insufficient conditions for realizability. In other words, the four Euclidean metric properties and relative-angle remain necessary but become a sufficient test only for positive semidefiniteness of all the principal  $3 \times 3$  submatrices [*sic*] in  $-V_N^T D V_N$ . Relative-angle inequality can be considered the ultimate

test only for realizability at each vertex  $x_k$  of each and every purported tetrahedron constituting a hyperdimensional body.

When  $N=5$  in particular, relative-angle inequality becomes the penultimate Euclidean metric requirement while nonnegativity of then unwieldy  $\det(\Theta^T \Theta)$  corresponds (by the *positive (semi)definite principal submatrices theorem* in §A.3.1.0.4) to the sixth and last Euclidean metric requirement. Together these six tests become necessary and sufficient, and so on.

Yet for all values of  $N$ , only assuming nonnegative  $d_{ij}$ , relative-angle matrix inequality in (1220) is necessary and sufficient to certify realizability; (§5.4.3.1)

$$\begin{array}{l} \text{Euclidean metric property 1 (§5.2).} \\ \text{Angle matrix inequality } \Omega \succeq 0 \text{ (1108).} \end{array} \Leftrightarrow \begin{array}{l} -V_N^T D V_N \succeq 0 \\ D \in \mathbb{S}_h^N \end{array} \Leftrightarrow D = \mathbf{D}(\Omega, d) \in \mathbb{EDM}^N \quad (1310)$$

Like matrix criteria (1028), (1052), and (1220), the relative-angle matrix inequality and nonnegativity property subsume all the Euclidean metric properties and further requirements.

### 5.14.3 Path not followed

As a means to test for realizability of four or more points, an intuitively appealing way to augment the four Euclidean metric properties is to recognize generalizations of the triangle inequality: In the case of cardinality  $N=4$  the three-dimensional analogue to triangle & distance is tetrahedron & facet-area, whereas in case  $N=5$  the four-dimensional analogue is polychoron & facet-volume, *ad infinitum*. For  $N$  points,  $N+1$  metric properties are required.

#### 5.14.3.1 $N = 4$

Each of the four facets of a general tetrahedron is a triangle and its relative interior. Suppose we identify each facet of the tetrahedron by its area-square:  $c_1, c_2, c_3, c_4$ . Then analogous to metric property 4, we may write a tight<sup>5.62</sup> area inequality for the facets

$$\sqrt{c_i} \leq \sqrt{c_j} + \sqrt{c_k} + \sqrt{c_\ell}, \quad i \neq j \neq k \neq \ell \in \{1, 2, 3, 4\} \quad (1311)$$

which is a generalized “triangle” inequality [254, §1.1] that follows from

$$\sqrt{c_i} = \sqrt{c_j} \cos \varphi_{ij} + \sqrt{c_k} \cos \varphi_{ik} + \sqrt{c_\ell} \cos \varphi_{i\ell} \quad (1312)$$

[271] [436, *Law of Cosines*] where  $\varphi_{ij}$  is the dihedral angle at the common edge between triangular facets  $i$  and  $j$ .

If  $D$  is the EDM corresponding to the whole tetrahedron, then area-square of the  $i^{\text{th}}$  triangular facet has a convenient formula in terms of  $D_i \in \mathbb{EDM}^{N-1}$  the EDM corresponding to that particular facet: From the *Cayley-Menger determinant*<sup>5.63</sup> for simplices, [436] [149] [186, §4] [96, §3.3] the  $i^{\text{th}}$  facet

<sup>5.62</sup>The upper bound is met when all angles in (1312) are simultaneously 0; that occurs, for example, if one point is relatively interior to the convex hull of the three remaining.

<sup>5.63</sup>whose foremost characteristic is: the determinant vanishes if and only if affine dimension does not equal penultimate cardinality; *id est*,  $\det \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} = 0 \Leftrightarrow r < N-1$  where  $D$  is any EDM (§5.7.3.0.1). Otherwise, the determinant is negative.

area-square for  $i \in \{1 \dots N\}$  is (§A.4.1)

$$c_i = \frac{-1}{2^{N-2}(N-2)!^2} \det \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D_i \end{bmatrix} \quad (1313)$$

$$= \frac{(-1)^N}{2^{N-2}(N-2)!^2} \det D_i (\mathbf{1}^T D_i^{-1} \mathbf{1}) \quad (1314)$$

$$= \frac{(-1)^N}{2^{N-2}(N-2)!^2} \mathbf{1}^T \text{cof}(D_i) \mathbf{1}^T \quad (1315)$$

where  $D_i$  is the  $i^{\text{th}}$  principal  $N-1 \times N-1$  submatrix<sup>5.64</sup> of  $D \in \mathbb{EDM}^N$ , and  $\text{cof}(D_i)$  is the  $N-1 \times N-1$  matrix of cofactors [368, §4] corresponding to  $D_i$ . The number of principal  $3 \times 3$  submatrices in  $D$  is, of course, equal to the number of triangular facets in the tetrahedron; four ( $N!/(3!(N-3)!)$ ) when  $N=4$ .

**5.14.3.1.1 Exercise.** *Sufficiency conditions for an EDM of four points.* Triangle inequality (property 4) and area inequality (1311) are conditions necessary for  $D$  to be an EDM. Prove their sufficiency in conjunction with the remaining three Euclidean metric properties. ▼

### 5.14.3.2 $N = 5$

Moving to the next level, we might encounter a Euclidean body called *polychoron*: a bounded polyhedron in four dimensions.<sup>5.65</sup> Our polychoron has five ( $N!/(4!(N-4)!)$ ) facets, each of them a general tetrahedron whose volume-square  $c_i$  is calculated using the same formula; (1313) where  $D$  is the EDM corresponding to the polychoron, and  $D_i$  is the EDM corresponding to the  $i^{\text{th}}$  facet (the principal  $4 \times 4$  submatrix of  $D \in \mathbb{EDM}^N$  corresponding to the  $i^{\text{th}}$  tetrahedron). The analogue to triangle & distance is now polychoron & facet-volume. We could then write another generalized “triangle” inequality like (1311) but in terms of facet volume; [441, §IV]

$$\sqrt{c_i} \leq \sqrt{c_j} + \sqrt{c_k} + \sqrt{c_\ell} + \sqrt{c_m}, \quad i \neq j \neq k \neq \ell \neq m \in \{1 \dots 5\} \quad (1316)$$

**5.14.3.2.1 Exercise.** *Sufficiency for an EDM of five points.*

For  $N=5$ , triangle (distance) inequality (§5.2), area inequality (1311), and volume inequality (1316) are conditions necessary for  $D$  to be an EDM. Prove their sufficiency. ▼

### 5.14.3.3 Volume of simplices

There is no known formula for the volume of a bounded general convex polyhedron expressed either by halfspace or vertex-description. [453, §2.1] [314, p.173] [261] [197] [198] Volume is a concept germane to  $\mathbb{R}^3$ ; in higher dimensions it is called *content*. Applying the *EDM assertion* (§5.9.1.0.4) and a result from [65, p.407], a general nonempty simplex (§2.12.3) in  $\mathbb{R}^{N-1}$  corresponding to an EDM  $D \in \mathbb{S}_h^N$  has content

$$\sqrt{c} = \text{content}(\mathcal{S}) \sqrt{\det(-V_N^T D V_N)} \quad (1317)$$

<sup>5.64</sup>Every principal submatrix of an EDM remains an EDM. [264, §4.1]

<sup>5.65</sup>The simplest polychoron is called a *pentatope* [436]; a regular simplex hence convex. (A *pentahedron* is a three-dimensional body having five vertices.)

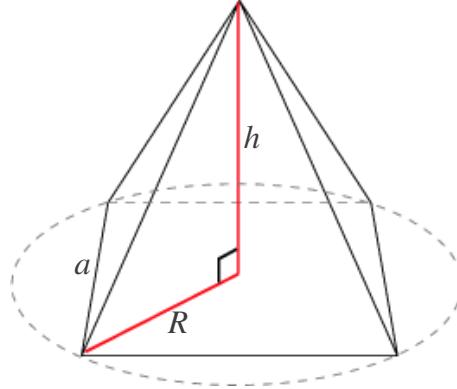


Figure 162: Length of one-dimensional face  $a$  equals height  $h=a=1$  of this convex nonsimplicial pyramid in  $\mathbb{R}^3$  with square base inscribed in a circle of radius  $R$  centered at the origin. [436, *Pyramid*]

where content-square of the unit simplex  $\mathcal{S} \subset \mathbb{R}^{N-1}$  is proportional to its Cayley-Menger determinant;

$$\text{content}(\mathcal{S})^2 = \frac{-1}{2^{N-1}(N-1)!^2} \det \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbf{D}([\mathbf{0} \ e_1 \ e_2 \ \dots \ e_{N-1}]) \end{bmatrix} \quad (1318)$$

where  $e_i \in \mathbb{R}^{N-1}$  and the EDM operator used is  $\mathbf{D}(X)$  (1033).

#### 5.14.3.3.1 Example. Pyramid.

A formula for volume of a pyramid is known:<sup>5.66</sup> it is  $\frac{1}{3}$  the product of its base area with its height. [250] The pyramid in Figure 162 has volume  $\frac{1}{3}$ . To find its volume using EDMs, we must first decompose the pyramid into simplicial parts. Slicing it in half along the plane containing the line segments corresponding to radius  $R$  and height  $h$  we find the vertices of one simplex,

$$X = \begin{bmatrix} 1/2 & 1/2 & -1/2 & 0 \\ 1/2 & -1/2 & -1/2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{n \times N} \quad (1319)$$

where  $N=n+1$  for any nonempty simplex in  $\mathbb{R}^n$ . The volume of this simplex is half that of the entire pyramid; *id est*,  $\sqrt{c}=\frac{1}{6}$  found by evaluating (1317).  $\square$

With that, we conclude digression of path.

#### 5.14.4 Affine dimension reduction in three dimensions

(confer §5.8.4) The determinant of any  $M \times M$  matrix is equal to the product of its  $M$  eigenvalues. [368, §5.1] When  $N=4$  and  $\det(\Theta^T \Theta)$  is 0, that means one or more eigenvalues of  $\Theta^T \Theta \in \mathbb{R}^{3 \times 3}$  are 0. The determinant will go to 0 whenever equality is attained on either side of (1027), (1308a), or (1308b), meaning that a tetrahedron has

---

<sup>5.66</sup>Pyramid volume is independent of the paramount vertex position as long as its height remains constant.

collapsed to a lower affine dimension; *id est*,  $r = \text{rank } \Theta^T \Theta = \text{rank } \Theta$  is reduced below  $N - 1$  exactly by the number of 0 eigenvalues ([§5.7.1.1](#)).

In solving completion problems of any size  $N$  where one or more entries of an EDM are unknown, therefore, dimension  $r$  of the affine hull required to contain the unknown points is potentially reduced by selecting distances to attain equality in [\(1027\)](#) or [\(1308a\)](#) or [\(1308b\)](#).

#### 5.14.4.1 *Exemplum redux*

We now apply the *fifth Euclidean metric property* to an earlier problem:

##### 5.14.4.1.1 Example. *Small completion problem, IV.* (confer [§5.9.3.0.1](#))

Returning again to [Example 5.3.0.0.2](#) that pertains to [Figure 141](#) where  $N=4$ , distance-square  $d_{14}$  is ascertainable from the fifth Euclidean metric property. Because all distances in [\(1025\)](#) are known except  $\sqrt{d_{14}}$ , then  $\cos \theta_{123} = 0$  and  $\theta_{324} = 0$  result from identity [\(1099\)](#). Applying [\(1027\)](#),

$$\begin{aligned} \cos(\theta_{123} + \theta_{324}) &\leq \cos \theta_{124} \leq \cos(\theta_{123} - \theta_{324}) \\ 0 &\leq \cos \theta_{124} \leq 0 \end{aligned} \tag{1320}$$

It follows again from [\(1099\)](#) that  $d_{14}$  can only be 2. As explained in this subsection, affine dimension  $r$  cannot exceed  $N - 2$  because equality is attained in [\(1320\)](#).  $\square$



# Chapter 6

## Cone of distance matrices

*For  $N > 3$ , the cone of EDMs is no longer a circular cone and the geometry becomes complicated...*

—Hayden, Wells, Liu, & Tarazaga, 1991 [211, §3]

In the subspace of symmetric matrices  $\mathbb{S}^N$ , we know that the convex cone of Euclidean distance matrices  $\mathbb{EDM}^N$  (the EDM cone) does not intersect the positive semidefinite cone  $\mathbb{S}_+^N$  (PSD cone) except at the origin, their only vertex; there can be no positive or negative semidefinite EDM. (1244) [264]

$$\mathbb{EDM}^N \cap \mathbb{S}_+^N = \mathbf{0} \quad (1321)$$

Even so, the two convex cones can be related. In §6.8.1 we prove the equality

$$\mathbb{EDM}^N = \mathbb{S}_h^N \cap \left( \mathbb{S}_c^{N\perp} - \mathbb{S}_+^N \right) \quad (1414)$$

a resemblance to EDM definition (1033) where

$$\mathbb{S}_h^N = \left\{ A \in \mathbb{S}^N \mid \delta(A) = \mathbf{0} \right\} \quad (67)$$

is the symmetric hollow subspace (§2.2.3) and where

$$\mathbb{S}_c^{N\perp} = \left\{ u \mathbf{1}^T + \mathbf{1} u^T \mid u \in \mathbb{R}^N \right\} \quad (2198)$$

is the orthogonal complement of the geometric center subspace (§E.7.2.0.2)

$$\mathbb{S}_c^N = \left\{ Y \in \mathbb{S}^N \mid Y \mathbf{1} = \mathbf{0} \right\} \quad (2196)$$

### 6.0.1 gravity

Equality (1414) is equally important as the known isomorphisms (1152) (1153) (1164) (1165) relating the EDM cone  $\mathbb{EDM}^N$  to positive semidefinite cone  $\mathbb{S}_+^{N-1}$  (§5.6.2.1) or to an  $N(N-1)/2$ -dimensional face of  $\mathbb{S}_+^N$  (§5.6.1.1).<sup>6.1</sup> But those isomorphisms have never led to this equality relating whole cones  $\mathbb{EDM}^N$  and  $\mathbb{S}_+^N$ .

Equality (1414) is not obvious from the various EDM definitions such as (1033) or (1337) because inclusion must be proved algebraically in order to establish equality;  $\mathbb{EDM}^N \supseteq \mathbb{S}_h^N \cap (\mathbb{S}_c^{N\perp} - \mathbb{S}_+^N)$ . We will instead prove (1414) using purely geometric methods.

### 6.0.2 highlight

In §6.8.1.7 we show: the Schoenberg criterion for discriminating Euclidean distance matrices

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} -V_N^T D V_N \in \mathbb{S}_+^{N-1} \\ D \in \mathbb{S}_h^N \end{cases} \quad (1052)$$

is a discretized membership relation (§2.13.4, dual generalized inequalities) between the EDM cone and its ordinary dual.

## 6.1 Defining EDM cone

We invoke a popular matrix criterion to illustrate correspondence between the EDM and PSD cones belonging to the ambient space of symmetric matrices:

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} -VDV \in \mathbb{S}_+^N \\ D \in \mathbb{S}_h^N \end{cases} \quad (1056)$$

where  $V \in \mathbb{S}^N$  is the geometric centering matrix (§B.4). The set of all EDMs of dimension  $N \times N$  forms a closed convex cone  $\mathbb{EDM}^N$  because any pair of EDMs satisfies the definition of a convex cone (178); *videlicet*, for each and every  $\zeta_1, \zeta_2 \geq 0$  (§A.3.1.0.2)

$$\begin{aligned} \zeta_1 V D_1 V + \zeta_2 V D_2 V &\succeq 0 &\Leftarrow& V D_1 V \succeq 0, \quad V D_2 V \succeq 0 \\ \zeta_1 D_1 + \zeta_2 D_2 &\in \mathbb{S}_h^N && D_1 \in \mathbb{S}_h^N, \quad D_2 \in \mathbb{S}_h^N \end{aligned} \quad (1322)$$

and convex cones are invariant to inverse linear transformation [343, p.22].

#### 6.1.0.0.1 Definition. Cone of Euclidean distance matrices.

In the subspace of symmetric matrices, the set of all Euclidean distance matrices forms a unique immutable pointed closed convex cone called the *EDM cone*: for  $N > 0$

$$\begin{aligned} \mathbb{EDM}^N &\triangleq \left\{ D \in \mathbb{S}_h^N \mid -VDV \in \mathbb{S}_+^N \right\} \\ &= \bigcap_{z \in \mathcal{N}(\mathbf{1}^T)} \left\{ D \in \mathbb{S}^N \mid \langle zz^T, -D \rangle \geq 0, \quad \delta(D) = \mathbf{0} \right\} \end{aligned} \quad (1323)$$

The EDM cone in isomorphic  $\mathbb{R}^{N(N+1)/2}$  [*sic*] is the intersection of an infinite number (when  $N > 2$ ) of halfspaces about the origin and a finite number of hyperplanes through the origin in vectorized variable  $D = [d_{ij}]$ . Hence  $\mathbb{EDM}^N$  is not full-dimensional with

---

<sup>6.1</sup>Because both positive semidefinite cones are frequently in play, dimension is explicit.

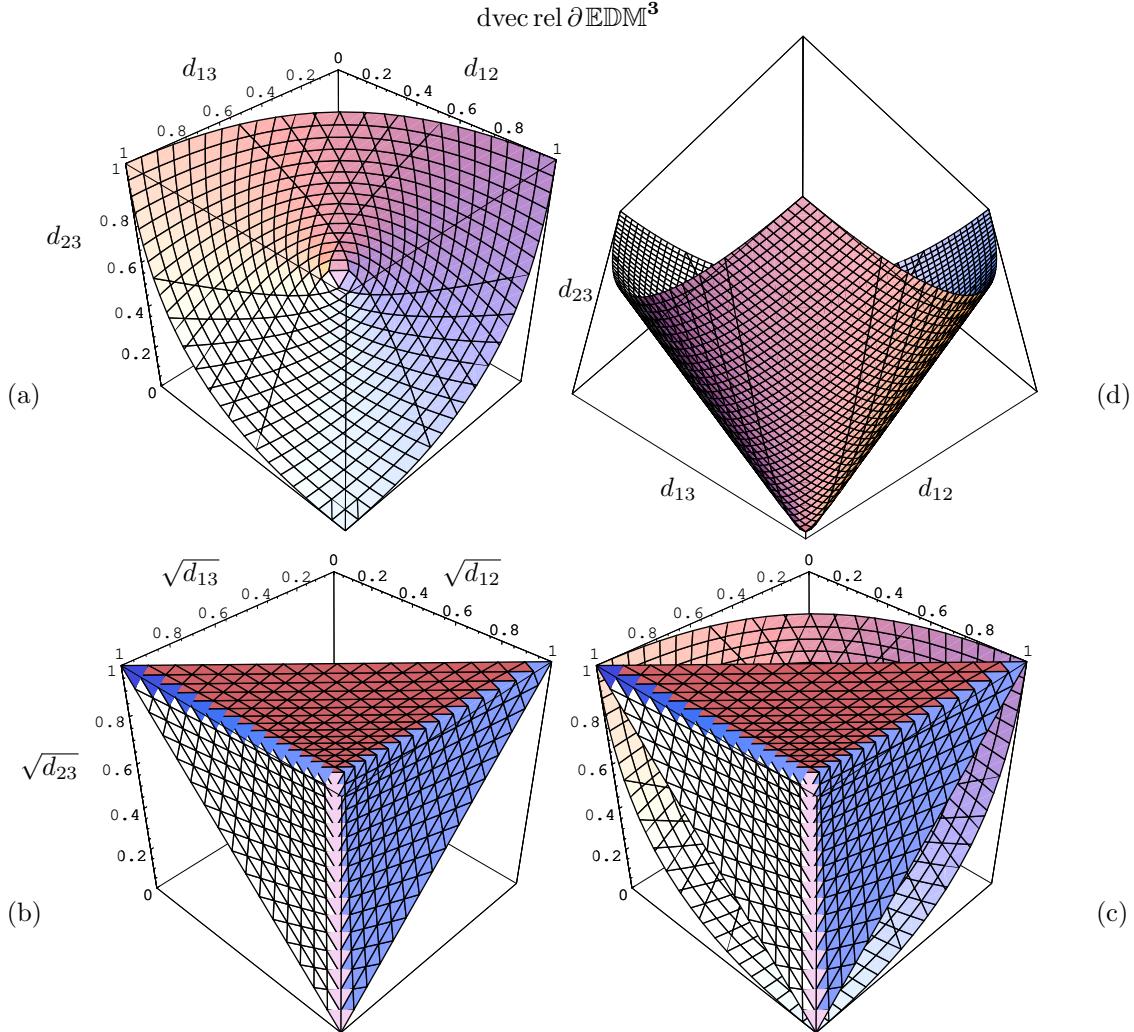


Figure 163: Relative boundary (tiled) of EDM cone  $\mathbb{EDM}^3$  drawn truncated in isometrically isomorphic subspace  $\mathbb{R}^3$ . **(a)** EDM cone drawn in usual distance-square coordinates  $d_{ij}$ . View is from interior toward origin. Unlike positive semidefinite cone, EDM cone is not selfdual; neither is it proper in ambient symmetric subspace (dual EDM cone for this example belongs to isomorphic  $\mathbb{R}^6$ ). **(b)** Drawn in its natural coordinates  $\sqrt{d_{ij}}$  (absolute distance), cone remains convex (*confer* §5.10); intersection of three halfspaces (1207) whose partial boundaries each contain origin. Cone geometry becomes nonconvex (nonpolyhedral) in higher dimension. (§6.3) **(c)** Two coordinate systems artificially superimposed. Coordinate transformation from  $d_{ij}$  to  $\sqrt{d_{ij}}$  appears a topological contraction. **(d)** Sitting on its vertex  $\mathbf{0}$ , pointed  $\mathbb{EDM}^3$  is a circular cone having axis of revolution  $\text{dvec}(-E) = \text{dvec}(\mathbf{1}\mathbf{1}^T - I)$  (1240) (74). (Rounded vertex is plot artifact.)

respect to  $\mathbb{S}^N$  because it is confined to the symmetric hollow subspace  $\mathbb{S}_h^N$ . The EDM cone relative interior comprises

$$\begin{aligned}\text{rel intr } \mathbb{EDM}^N &= \bigcap_{z \in \mathcal{N}(\mathbf{1}^T)} \left\{ D \in \mathbb{S}^N \mid \langle zz^T, -D \rangle > 0, \quad \delta(D) = \mathbf{0} \right\} \\ &= \left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) = N-1 \right\}\end{aligned}\tag{1324}$$

while its relative boundary comprises

$$\begin{aligned}\text{rel } \partial \mathbb{EDM}^N &= \left\{ D \in \mathbb{EDM}^N \mid \langle zz^T, -D \rangle = 0 \text{ for some } z \in \mathcal{N}(\mathbf{1}^T) \right\} \\ &= \left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) < N-1 \right\}\end{aligned}\tag{1325}$$

△

This cone is more easily visualized in the isomorphic vector subspace  $\mathbb{R}^{N(N-1)/2}$  corresponding to  $\mathbb{S}_h^N$ :

In the case  $N=1$  point, the EDM cone is the origin in  $\mathbb{R}^0$ .

In the case  $N=2$ , the EDM cone is the nonnegative real line in  $\mathbb{R}$ ; a halfline in a subspace of the realization in Figure 167.

The EDM cone in the case  $N=3$  is a circular cone in  $\mathbb{R}^3$  illustrated in Figure 163(a)(d); rather, the set of all matrices

$$D = \begin{bmatrix} 0 & d_{12} & d_{13} \\ d_{12} & 0 & d_{23} \\ d_{13} & d_{23} & 0 \end{bmatrix} \in \mathbb{EDM}^3\tag{1326}$$

makes a circular cone in this dimension. In this case, the first four Euclidean metric properties are necessary and sufficient tests to certify realizability of triangles; (1302). Thus triangle inequality property 4 describes three halfspaces (1207) whose intersection makes a polyhedral cone in  $\mathbb{R}^3$  of realizable  $\sqrt{d_{ij}}$  (absolute distance); an isomorphic subspace representation of the set of all EDMs  $D$  in the natural coordinates

$$\sqrt[3]{D} \triangleq \begin{bmatrix} 0 & \sqrt{d_{12}} & \sqrt{d_{13}} \\ \sqrt{d_{12}} & 0 & \sqrt{d_{23}} \\ \sqrt{d_{13}} & \sqrt{d_{23}} & 0 \end{bmatrix}\tag{1327}$$

illustrated in Figure 163b.

## 6.2 Polyhedral bounds

The convex cone of EDMs is nonpolyhedral in  $d_{ij}$  for  $N > 2$ ; e.g., Figure 163a. Still we found necessary and sufficient bounding polyhedral relations consistent with EDM cones for cardinality  $N = 1, 2, 3, 4$ :

- $N=3$ . Transforming distance-square coordinates  $d_{ij}$  by taking their positive square root provides the polyhedral cone in Figure 163b; polyhedral because an intersection of three halfspaces in natural coordinates  $\sqrt{d_{ij}}$  is provided by triangle inequalities (1207). This polyhedral cone implicitly encompasses necessary and sufficient metric properties: nonnegativity, selfdistance, symmetry, and triangle inequality.
- $N=4$ . Relative-angle inequality (1308) together with four Euclidean metric properties are necessary and sufficient tests for realizability of tetrahedra. (1309) Albeit relative angles  $\theta_{ikj}$  (1099) are nonlinear functions of the  $d_{ij}$ , relative-angle inequality

provides a regular tetrahedron in  $\mathbb{R}^3$  [sic] (Figure 161) bounding angles  $\theta_{ikj}$  at vertex  $x_k$  consistently with  $\text{EDM}^4$ .<sup>6.2</sup>

Yet were we to employ the procedure outlined in §5.14.3 for making generalized triangle inequalities, then we would find all the necessary and sufficient  $d_{ij}$ -transformations for generating bounding polyhedra consistent with EDMs of any higher dimension ( $N > 3$ ).

### 6.3 $\sqrt{\text{EDM}}$ cone is not convex

For some applications, like a molecular conformation problem (Figure 5, Figure 152) or multidimensional scaling [115] [395], absolute distance  $\sqrt{d_{ij}}$  is the preferred variable. Taking square root of the entries in all EDMs  $D$  of dimension  $N$ , we get another cone but not a convex cone when  $N > 3$  (Figure 163b): [97, §4.5.2]

$$\sqrt{\text{EDM}}^N \triangleq \{\sqrt[3]{D} \mid D \in \text{EDM}^N\} \quad (1328)$$

where  $\sqrt[3]{D}$  is defined like (1327). It is a cone because any cone is completely constituted by rays emanating from the origin: (§2.7) Any given ray  $\{\zeta \Gamma \in \mathbb{R}^{N(N-1)/2} \mid \zeta \geq 0\}$  remains a ray  $\{\sqrt[3]{\zeta \Gamma} \in \mathbb{R}^{N(N-1)/2} \mid \zeta \geq 0\}$  under entrywise square root. It is already established that

$$D \in \text{EDM}^N \Rightarrow \sqrt[3]{D} \in \text{EDM}^N \quad (1239)$$

But because of how  $\sqrt{\text{EDM}}^N$  is defined, it is obvious that (confer §5.10)

$$D \in \text{EDM}^N \Leftrightarrow \sqrt[3]{D} \in \sqrt{\text{EDM}}^N \quad (1329)$$

Were  $\sqrt{\text{EDM}}^N$  convex, then given  $\sqrt[3]{D_1}, \sqrt[3]{D_2} \in \sqrt{\text{EDM}}^N$  we would expect their conic combination  $\sqrt[3]{D_1} + \sqrt[3]{D_2}$  to be a member of  $\sqrt{\text{EDM}}^N$ . That is easily proven false by counterexample via (1329), for then  $(\sqrt[3]{D_1} + \sqrt[3]{D_2}) \circ (\sqrt[3]{D_1} + \sqrt[3]{D_2})$  would need to be a member of  $\text{EDM}^N$ .

Notwithstanding,

$$\sqrt{\text{EDM}}^N \subseteq \text{EDM}^N \quad (1330)$$

by (1239) (Figure 163), and we learn how to transform a nonconvex *proximity problem* in the natural coordinates  $\sqrt{d_{ij}}$  to a convex optimization in §7.2.1.

### 6.4 EDM definition in $11^T$

Any EDM  $D$  corresponding to affine dimension  $r$  has representation

$$\mathbf{D}(V_{\mathcal{X}}) \triangleq \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) \mathbf{1}^T + \mathbf{1} \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T)^T - 2V_{\mathcal{X}} V_{\mathcal{X}}^T \in \text{EDM}^N \quad (1331)$$

where  $\mathcal{R}(V_{\mathcal{X}} \in \mathbb{R}^{N \times r}) \subseteq \mathcal{N}(\mathbf{1}^T) = \mathbf{1}^\perp$

$$V_{\mathcal{X}}^T V_{\mathcal{X}} = \delta^2(V_{\mathcal{X}}^T V_{\mathcal{X}}) \quad \text{and} \quad V_{\mathcal{X}} \text{ is full-rank with orthogonal columns.} \quad (1332)$$

Equation (1331) is simply the standard EDM definition (1033) with a centered list  $X$  as in (1120); Gram matrix  $X^T X$  has been replaced with the subcompact singular value decomposition (§A.6.2)<sup>6.3</sup>

$$V_{\mathcal{X}} V_{\mathcal{X}}^T \equiv V^T X^T X V \in \mathbb{S}_c^N \cap \mathbb{S}_+^N \quad (1333)$$

<sup>6.2</sup>Still, property-4 triangle inequalities (1207) corresponding to each principal  $3 \times 3$  submatrix of  $-V_{\mathcal{X}}^T D V_{\mathcal{X}}$  demand that the corresponding  $\sqrt{d_{ij}}$  belong to a polyhedral cone like that in Figure 163b.

<sup>6.3</sup>Subcompact SVD:  $V_{\mathcal{X}} V_{\mathcal{X}}^T \triangleq Q \sqrt{\Sigma} \sqrt{\Sigma} Q^T \equiv V^T X^T X V$ . So  $V_{\mathcal{X}}^T$  is not necessarily  $X V$  (§5.5.1.0.1), although affine dimension  $r = \text{rank}(V_{\mathcal{X}}^T) = \text{rank}(X V)$ . (1177)

This means: inner product  $V_{\mathcal{X}}^T V_{\mathcal{X}}$  is an  $r \times r$  diagonal matrix  $\Sigma$  of nonzero singular values.

Vector  $\delta(V_{\mathcal{X}} V_{\mathcal{X}}^T)$  may be decomposed into complementary parts by projecting it on orthogonal subspaces  $\mathbf{1}^\perp$  and  $\mathcal{R}(\mathbf{1})$ : namely,

$$P_{\mathbf{1}^\perp}(\delta(V_{\mathcal{X}} V_{\mathcal{X}}^T)) = V \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) \quad (1334)$$

$$P_{\mathbf{1}}(\delta(V_{\mathcal{X}} V_{\mathcal{X}}^T)) = \frac{1}{N} \mathbf{1} \mathbf{1}^T \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) \quad (1335)$$

Of course

$$\delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) = V \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) + \frac{1}{N} \mathbf{1} \mathbf{1}^T \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) \quad (1336)$$

by (1055). Substituting this into EDM definition (1331), we get the Hayden, Wells, Liu, & Tarazaga EDM formula [211, §2]

$$\mathbf{D}(V_{\mathcal{X}}, y) \triangleq y \mathbf{1}^T + \mathbf{1} y^T + \frac{\lambda}{N} \mathbf{1} \mathbf{1}^T - 2 V_{\mathcal{X}} V_{\mathcal{X}}^T \in \mathbb{EDM}^N \quad (1337)$$

where

$$\lambda \triangleq 2 \|V_{\mathcal{X}}\|_F^2 = \mathbf{1}^T \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) 2 \quad \text{and} \quad y \triangleq \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) - \frac{\lambda}{2N} \mathbf{1} = V \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) \quad (1338)$$

and  $y = \mathbf{0}$  if and only if  $\mathbf{1}$  is an eigenvector of EDM  $D$ . Scalar  $\lambda$  becomes an eigenvalue when corresponding eigenvector  $\mathbf{1}$  exists.<sup>6.4</sup>

Then the particular dyad sum from (1337)

$$y \mathbf{1}^T + \mathbf{1} y^T + \frac{\lambda}{N} \mathbf{1} \mathbf{1}^T \in \mathbb{S}_c^{N \perp} \quad (1339)$$

must belong to the orthogonal complement of the geometric center subspace (p.593), whereas  $V_{\mathcal{X}} V_{\mathcal{X}}^T \in \mathbb{S}_c^N \cap \mathbb{S}_+^N$  (1333) belongs to the positive semidefinite cone in the geometric center subspace.

**Proof.** We validate eigenvector  $\mathbf{1}$  and eigenvalue  $\lambda$ .

( $\Rightarrow$ ) Suppose  $\mathbf{1}$  is an eigenvector of EDM  $D$ . Then because

$$V_{\mathcal{X}}^T \mathbf{1} = \mathbf{0} \quad (1340)$$

it follows

$$\begin{aligned} D \mathbf{1} &= \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) \mathbf{1}^T \mathbf{1} + \mathbf{1} \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T)^T \mathbf{1} = N \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) + \|V_{\mathcal{X}}\|_F^2 \mathbf{1} \\ &\Rightarrow \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) \propto \mathbf{1} \end{aligned} \quad (1341)$$

For some  $\kappa \in \mathbb{R}_+$

$$\delta(V_{\mathcal{X}} V_{\mathcal{X}}^T)^T \mathbf{1} = N \kappa = \text{tr}(V_{\mathcal{X}}^T V_{\mathcal{X}}) = \|V_{\mathcal{X}}\|_F^2 \Rightarrow \delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) = \frac{1}{N} \|V_{\mathcal{X}}\|_F^2 \mathbf{1} \quad (1342)$$

so  $y = \mathbf{0}$ .

( $\Leftarrow$ ) Now suppose  $\delta(V_{\mathcal{X}} V_{\mathcal{X}}^T) = \frac{\lambda}{2N} \mathbf{1}$ ; id est,  $y = \mathbf{0}$ . Then

$$D = \frac{\lambda}{N} \mathbf{1} \mathbf{1}^T - 2 V_{\mathcal{X}} V_{\mathcal{X}}^T \in \mathbb{EDM}^N \quad (1343)$$

$\mathbf{1}$  is an eigenvector with corresponding eigenvalue  $\lambda$ . ♦

<sup>6.4</sup> e.g., when  $X = I$  in EDM definition (1033).

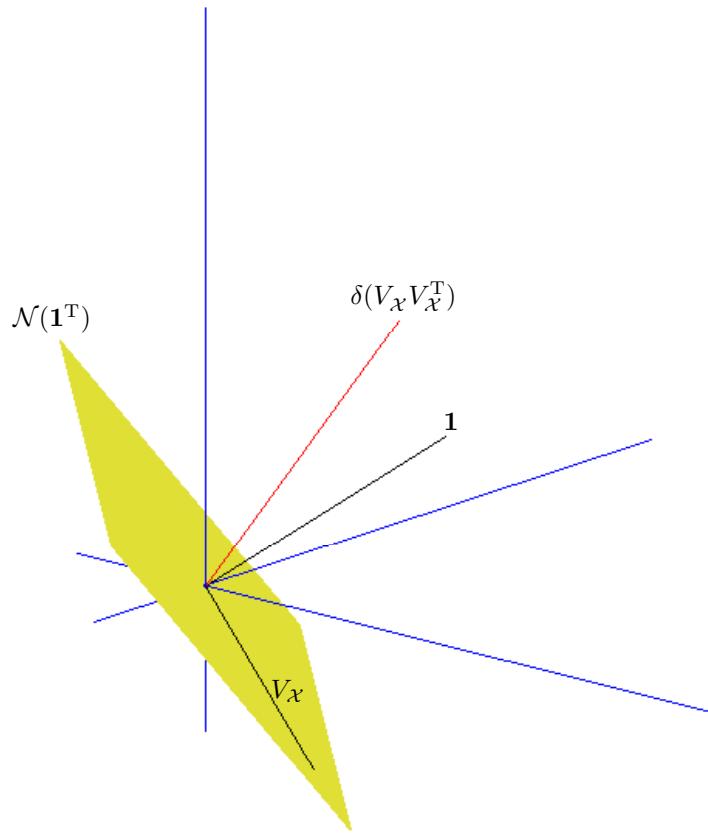


Figure 164: Example of  $V_x$  selection to make an EDM corresponding to cardinality  $N=3$  and affine dimension  $r=1$ ;  $V_x$  is a vector in nullspace  $\mathcal{N}(\mathbf{1}^T) \subset \mathbb{R}^3$ . Nullspace of  $\mathbf{1}^T$  is hyperplane in  $\mathbb{R}^3$  (drawn truncated) having normal  $\mathbf{1}$ . Vector  $\delta(V_x V_x^T)$  may or may not be in plane spanned by  $\{\mathbf{1}, V_x\}$ , but belongs to nonnegative orthant which is strictly supported by  $\mathcal{N}(\mathbf{1}^T)$ .

### 6.4.1 Range of EDM $D$

From §B.1.1 pertaining to linear independence of dyad sums: If the transpose halves of all the dyads in the sum (1331)<sup>6.5</sup> make a linearly independent set, then the nontranspose halves constitute a basis for the range of EDM  $D$ . Saying this mathematically: For  $D \in \mathbb{EDM}^N$

$$\begin{aligned}\mathcal{R}(D) = \mathcal{R}([\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) \quad \mathbf{1} \quad V_{\mathcal{X}}]) &\Leftarrow \text{rank}([\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) \quad \mathbf{1} \quad V_{\mathcal{X}}]) = 2 + r \\ \mathcal{R}(D) = \mathcal{R}([\mathbf{1} \quad V_{\mathcal{X}}]) &\Leftarrow \text{otherwise}\end{aligned}\quad (1344)$$

**6.4.1.0.1 Proof.** We need that condition under which the rank equality is satisfied: We know  $\mathcal{R}(V_{\mathcal{X}}) \perp \mathbf{1}$ , but what is the relative geometric orientation of  $\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T)$ ?  $\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) \succeq 0$  because  $V_{\mathcal{X}}V_{\mathcal{X}}^T \succeq 0$ , and  $\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) \propto \mathbf{1}$  remains possible (1341); this means  $\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) \notin \mathcal{N}(\mathbf{1}^T)$  simply because it has no negative entries. (Figure 164) If the projection of  $\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T)$  on  $\mathcal{N}(\mathbf{1}^T)$  does not belong to  $\mathcal{R}(V_{\mathcal{X}})$ , then that is a necessary and sufficient condition for linear independence (l.i.) of  $\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T)$  with respect to  $\mathcal{R}([\mathbf{1} \quad V_{\mathcal{X}}])$ ; *id est*,

$$\begin{aligned}V\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) &\neq V_{\mathcal{X}}a \quad \text{for any } a \in \mathbb{R}^r \\ (I - \frac{1}{N}\mathbf{1}\mathbf{1}^T)\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) &\neq V_{\mathcal{X}}a \\ \delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) - \frac{1}{N}\|V_{\mathcal{X}}\|_F^2\mathbf{1} &\neq V_{\mathcal{X}}a \\ \delta(V_{\mathcal{X}}V_{\mathcal{X}}^T) - \frac{\lambda}{2N}\mathbf{1} = y &\neq V_{\mathcal{X}}a \Leftrightarrow \{\mathbf{1}, \delta(V_{\mathcal{X}}V_{\mathcal{X}}^T), V_{\mathcal{X}}\} \text{ is l.i.}\end{aligned}\quad (1345)$$

When this condition is violated (when (1338)  $y = V_{\mathcal{X}}a_p$  for some particular  $a \in \mathbb{R}^r$ ), on the other hand, then from (1337) we have

$$\begin{aligned}\mathcal{R}(D = y\mathbf{1}^T + \mathbf{1}y^T + \frac{\lambda}{N}\mathbf{1}\mathbf{1}^T - 2V_{\mathcal{X}}V_{\mathcal{X}}^T) &= \mathcal{R}((V_{\mathcal{X}}a_p + \frac{\lambda}{N}\mathbf{1})\mathbf{1}^T + (\mathbf{1}a_p^T - 2V_{\mathcal{X}})V_{\mathcal{X}}^T) \\ &= \mathcal{R}([V_{\mathcal{X}}a_p + \frac{\lambda}{N}\mathbf{1} \quad \mathbf{1}a_p^T - 2V_{\mathcal{X}}]) \\ &= \mathcal{R}([\mathbf{1} \quad V_{\mathcal{X}}])\end{aligned}\quad (1346)$$

An example of such a violation is (1343) where, in particular,  $a_p = \mathbf{0}$ . ♦

Then a statement parallel to (1344) is, for  $D \in \mathbb{EDM}^N$  (Theorem 5.7.3.0.1)

$$\begin{aligned}\text{rank}(D) = r + 2 &\Leftrightarrow y \notin \mathcal{R}(V_{\mathcal{X}}) \quad (\Leftrightarrow \mathbf{1}^T D^\dagger \mathbf{1} = 0) \\ \text{rank}(D) = r + 1 &\Leftrightarrow y \in \mathcal{R}(V_{\mathcal{X}}) \quad (\Leftrightarrow \mathbf{1}^T D^\dagger \mathbf{1} \neq 0)\end{aligned}\quad (1347)$$

### 6.4.2 Boundary constituents of EDM cone

Expression (1331) has utility in forming the set of all EDMs corresponding to affine dimension  $r$ :

$$\begin{aligned}&\left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) = r \right\} \\ &= \left\{ \mathbf{D}(V_{\mathcal{X}}) \mid V_{\mathcal{X}} \in \mathbb{R}^{N \times r}, \text{rank } V_{\mathcal{X}} = r, V_{\mathcal{X}}^T V_{\mathcal{X}} = \delta^2(V_{\mathcal{X}}^T V_{\mathcal{X}}), \mathcal{R}(V_{\mathcal{X}}) \subseteq \mathcal{N}(\mathbf{1}^T) \right\}\end{aligned}\quad (1348)$$

whereas  $\{D \in \mathbb{EDM}^N \mid \text{rank}(VDV) \leq r\}$  is the closure of this same set;

$$\overline{\left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) \leq r \right\}} = \overline{\left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) = r \right\}}\quad (1349)$$

<sup>6.5</sup>Identifying columns  $V_{\mathcal{X}} \triangleq [v_1 \cdots v_r]$ , then  $V_{\mathcal{X}}V_{\mathcal{X}}^T = \sum_i v_i v_i^T$  is also a sum of dyads.

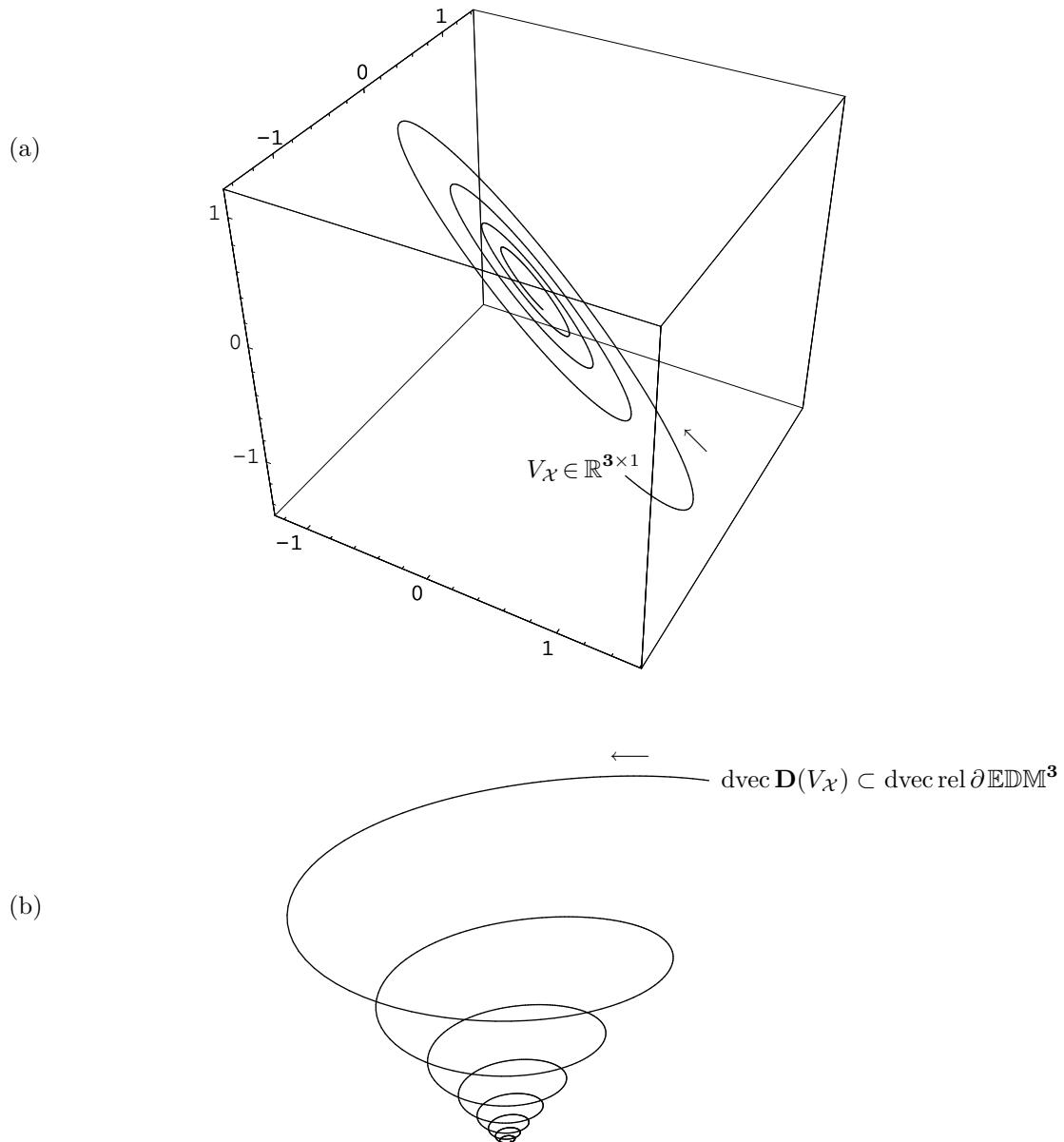


Figure 165: (a) Vector  $V_X$  from Figure 164 spirals in  $\mathcal{N}(\mathbf{1}^T) \subset \mathbb{R}^3$  decaying toward origin. (Spiral is two-dimensional in vector space  $\mathbb{R}^3$ .) (b) Corresponding trajectory  $\mathbf{D}(V_X)$  on EDM cone relative boundary creates a vortex also decaying toward origin. There are two complete orbits on EDM cone boundary about axis of revolution for every single revolution of  $V_X$  about origin. (Vortex is three-dimensional in isometrically isomorphic  $\mathbb{R}^3$ .)

For example,

$$\begin{aligned}\text{rel } \partial \mathbb{EDM}^N &= \left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) < N-1 \right\} \\ &= \bigcup_{r=0}^{N-2} \left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) = r \right\}\end{aligned}\quad (1350)$$

None of these are necessarily convex sets, although

$$\begin{aligned}\mathbb{EDM}^N &= \bigcup_{r=0}^{N-1} \left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) = r \right\} \\ &= \left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) = N-1 \right\} \\ \text{rel intr } \mathbb{EDM}^N &= \left\{ D \in \mathbb{EDM}^N \mid \text{rank}(VDV) = N-1 \right\}\end{aligned}\quad (1351)$$

are pointed convex cones.

When cardinality  $N = 3$  and affine dimension  $r = 2$ , for example, the relative interior  $\text{rel intr } \mathbb{EDM}^3$  is realized via (1348). (§6.5)

When  $N = 3$  and  $r = 1$ , the relative boundary of the EDM cone  $\text{dvec rel } \partial \mathbb{EDM}^3$  is realized in isomorphic  $\mathbb{R}^3$  as in Figure 163d. This figure could be constructed via (1349) by spiraling vector  $V_X$  tightly about the origin in  $\mathcal{N}(\mathbf{1}^T)$ ; as can be imagined with aid of Figure 164. Vectors close to the origin in  $\mathcal{N}(\mathbf{1}^T)$  are correspondingly close to the origin in  $\mathbb{EDM}^N$ . As vector  $V_X$  orbits the origin in  $\mathcal{N}(\mathbf{1}^T)$ , the corresponding EDM orbits the axis of revolution while remaining on the boundary of the circular cone  $\text{dvec rel } \partial \mathbb{EDM}^3$ . (Figure 165)

### 6.4.3 Faces of EDM cone

Like the positive semidefinite cone, EDM cone faces are EDM cones.

#### 6.4.3.0.1 Exercise. Isomorphic faces.

Prove that in high cardinality  $N$ , any set of EDMs made via (1348) or (1349) with particular affine dimension  $r$  is isomorphic with any set admitting the same affine dimension but made in lower cardinality. ▼

#### 6.4.3.1 smallest face that contains an EDM

Now suppose we are given a particular EDM  $\mathbf{D}(V_{X_p}) \in \mathbb{EDM}^N$  corresponding to affine dimension  $r$  and parametrized by  $V_{X_p}$  in (1331). The EDM cone's smallest face that contains  $\mathbf{D}(V_{X_p})$  is

$$\begin{aligned}\mathcal{F}\left(\mathbb{EDM}^N \ni \mathbf{D}(V_{X_p})\right) \\ &= \overline{\left\{ \mathbf{D}(V_X) \mid V_X \in \mathbb{R}^{N \times r}, \text{rank } V_X = r, V_X^T V_X = \delta^2(V_X^T V_X), \mathcal{R}(V_X) \subseteq \mathcal{R}(V_{X_p}) \right\}} \\ &\simeq \mathbb{EDM}^{r+1}\end{aligned}\quad (1352)$$

which is isomorphic<sup>6.6</sup> with convex cone  $\mathbb{EDM}^{r+1}$ , hence of dimension

$$\dim \mathcal{F}\left(\mathbb{EDM}^N \ni \mathbf{D}(V_{X_p})\right) = r(r+1)/2 \quad (1353)$$

---

<sup>6.6</sup>The fact that the smallest face is isomorphic with another EDM cone (perhaps smaller than  $\mathbb{EDM}^N$ ) is implicit in [211, §2].

in isomorphic  $\mathbb{R}^{N(N-1)/2}$ . Not all dimensions are represented; *e.g.*, the EDM cone has no two-dimensional faces.

When cardinality  $N=4$  and affine dimension  $r=2$  so that  $\mathcal{R}(V_{\mathcal{X}_p})$  is any two-dimensional subspace of three-dimensional  $\mathcal{N}(\mathbf{1}^T)$  in  $\mathbb{R}^4$ , for example, then the corresponding face of  $\mathbb{EDM}^4$  is isometrically isomorphic with: (1349)

$$\mathbb{EDM}^3 = \{D \in \mathbb{EDM}^3 \mid \text{rank}(VDV) \leq 2\} \simeq \mathcal{F}(\mathbb{EDM}^4 \ni \mathbf{D}(V_{\mathcal{X}_p})) \quad (1354)$$

Each two-dimensional subspace of  $\mathcal{N}(\mathbf{1}^T)$  corresponds to another three-dimensional face.

Because each and every principal submatrix of an EDM in  $\mathbb{EDM}^N$  (§5.14.3) is another EDM [264, §4.1], for example, then each principal submatrix belongs to a particular face of  $\mathbb{EDM}^N$ .

#### 6.4.3.2 extreme directions of EDM cone

In particular, extreme directions (§2.8.1) of  $\mathbb{EDM}^N$  correspond to affine dimension  $r=1$  and are simply represented: for any particular cardinality  $N \geq 2$  (§2.8.2) and each and every nonzero vector  $z$  in  $\mathcal{N}(\mathbf{1}^T)$

$$\begin{aligned} \Gamma &\triangleq (z \circ z)\mathbf{1}^T + \mathbf{1}(z \circ z)^T - 2zz^T \in \mathbb{EDM}^N \\ &= \delta(zz^T)\mathbf{1}^T + \mathbf{1}\delta(zz^T)^T - 2zz^T \end{aligned} \quad (1355)$$

is an extreme direction corresponding to a one-dimensional face of the EDM cone  $\mathbb{EDM}^N$  that is a ray in isomorphic subspace  $\mathbb{R}^{N(N-1)/2}$ .

Proving this would exercise the fundamental definition (189) of extreme direction. Here is a sketch: Any EDM may be represented

$$\mathbf{D}(V_{\mathcal{X}}) = \delta(V_{\mathcal{X}}V_{\mathcal{X}}^T)\mathbf{1}^T + \mathbf{1}\delta(V_{\mathcal{X}}V_{\mathcal{X}}^T)^T - 2V_{\mathcal{X}}V_{\mathcal{X}}^T \in \mathbb{EDM}^N \quad (1331)$$

where matrix  $V_{\mathcal{X}}$  (1332) has orthogonal columns. For the same reason (1621) that  $zz^T$  is an extreme direction of the positive semidefinite cone (§2.9.2.7) for any particular nonzero vector  $z$ , there is no conic combination of distinct EDMs (each conically independent of  $\Gamma$  (§2.10)) equal to  $\Gamma$ . ■

##### 6.4.3.2.1 Example. Biorthogonal expansion of an EDM. (confer §2.13.8.1.1)

When matrix  $D$  belongs to the EDM cone, nonnegative coordinates for biorthogonal expansion are the eigenvalues  $\lambda \in \mathbb{R}^N$  of  $-VDV\frac{1}{2}$ : For any  $D \in \mathbb{S}_h^N$  it holds

$$D = \delta(-VDV\frac{1}{2})\mathbf{1}^T + \mathbf{1}\delta(-VDV\frac{1}{2})^T - 2(-VDV\frac{1}{2}) \quad (1144)$$

By diagonalization  $-VDV\frac{1}{2} \triangleq Q\Lambda Q^T \in \mathbb{S}_c^N$  (§A.5.1) we may write

$$\begin{aligned} D &= \delta\left(\sum_{i=1}^N \lambda_i q_i q_i^T\right)\mathbf{1}^T + \mathbf{1}\delta\left(\sum_{i=1}^N \lambda_i q_i q_i^T\right)^T - 2\sum_{i=1}^N \lambda_i q_i q_i^T \\ &= \sum_{i=1}^N \lambda_i (\delta(q_i q_i^T)\mathbf{1}^T + \mathbf{1}\delta(q_i q_i^T)^T - 2q_i q_i^T) \end{aligned} \quad (1356)$$

where  $q_i$  is the  $i^{\text{th}}$  eigenvector of  $-VDV\frac{1}{2}$  arranged columnar in orthogonal matrix

$$Q = [q_1 \ q_2 \ \cdots \ q_N] \in \mathbb{R}^{N \times N} \quad (406)$$

and where  $\{\delta(q_i q_i^T) \mathbf{1}^T + \mathbf{1} \delta(q_i q_i^T)^T - 2q_i q_i^T, i=1 \dots N\}$  are extreme directions of some pointed polyhedral cone  $\mathcal{K} \subset \mathbb{S}_h^N$  and extreme directions of  $\mathbb{EDM}^N$ . Invertibility of (1356)

$$\begin{aligned} -VDV^{\frac{1}{2}} &= -V \sum_{i=1}^N \lambda_i (\delta(q_i q_i^T) \mathbf{1}^T + \mathbf{1} \delta(q_i q_i^T)^T - 2q_i q_i^T) V^{\frac{1}{2}} \\ &= \sum_{i=1}^N \lambda_i q_i q_i^T \end{aligned} \quad (1357)$$

implies linear independence of those extreme directions. Then biorthogonal expansion is expressed

$$\text{dvec } D = YY^\dagger \text{dvec } D = Y \lambda(-VDV^{\frac{1}{2}}) \quad (1358)$$

where

$$Y \triangleq [\text{dvec}(\delta(q_i q_i^T) \mathbf{1}^T + \mathbf{1} \delta(q_i q_i^T)^T - 2q_i q_i^T), i=1 \dots N] \in \mathbb{R}^{N(N-1)/2 \times N} \quad (1359)$$

When  $D$  belongs to the EDM cone in the subspace of symmetric hollow matrices, unique coordinates  $Y^\dagger \text{dvec } D$  for this biorthogonal expansion must be the nonnegative eigenvalues  $\lambda$  of  $-VDV^{\frac{1}{2}}$ . This means  $D$  simultaneously belongs to the EDM cone and to the pointed polyhedral cone  $\text{dvec } \mathcal{K} = \text{cone}(Y)$ .  $\square$

#### 6.4.3.3 open question

Result (1353) is analogous to that for the positive semidefinite cone (226), although the question remains open whether all faces of  $\mathbb{EDM}^N$  (whose dimension is less than dimension of the cone) are exposed like they are for the positive semidefinite cone.<sup>6.7</sup> (§2.9.2.4) [386]

## 6.5 Correspondence to PSD cone $\mathbb{S}_+^{N-1}$

Hayden, Wells, Liu, & Tarazaga [211, §2] assert one-to-one correspondence of EDMs with positive semidefinite matrices in the symmetric subspace. Because  $\text{rank}(VDV) \leq N-1$  (§5.7.1.1), that PSD cone corresponding to the EDM cone can only be  $\mathbb{S}_+^{N-1}$ . [9, §18.2.1] To clearly demonstrate this correspondence, we invoke inner-product form EDM definition

$$\begin{aligned} \mathbf{D}(\Phi) &= \begin{bmatrix} 0 \\ \delta(\Phi) \end{bmatrix} \mathbf{1}^T + \mathbf{1} \begin{bmatrix} 0 & \delta(\Phi)^T \end{bmatrix} - 2 \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \Phi \end{bmatrix} \in \mathbb{EDM}^N \\ &\Leftrightarrow \\ \Phi &\succeq 0 \end{aligned} \quad (1162)$$

Then the EDM cone may be expressed

$$\mathbb{EDM}^N = \left\{ \mathbf{D}(\Phi) \mid \Phi \in \mathbb{S}_+^{N-1} \right\} \quad (1360)$$

Hayden & Wells' assertion can therefore be equivalently stated in terms of an inner-product form EDM operator

$$\mathbf{D}(\mathbb{S}_+^{N-1}) = \mathbb{EDM}^N \quad (1164)$$

$$\mathbf{V}_{\mathcal{N}}(\mathbb{EDM}^N) = \mathbb{S}_+^{N-1} \quad (1165)$$

identity (1165) holding because  $\mathcal{R}(V_{\mathcal{N}}) = \mathcal{N}(\mathbf{1}^T)$  (1040), linear functions  $\mathbf{D}(\Phi)$  and  $\mathbf{V}_{\mathcal{N}}(D) = -V_{\mathcal{N}}^T D V_{\mathcal{N}}$  (§5.6.2.1) being mutually inverse.

In terms of affine dimension  $r$ , Hayden & Wells claim particular correspondence between PSD and EDM cones:

---

<sup>6.7</sup>Elementary example of a face not exposed is given by the closed convex set in Figure 35 and in Figure 45.

$r = N - 1$ : Symmetric hollow matrices  $-D$  positive definite on  $\mathcal{N}(\mathbf{1}^T)$  correspond to points relatively interior to the EDM cone.

$r < N - 1$ : Symmetric hollow matrices  $-D$  positive semidefinite on  $\mathcal{N}(\mathbf{1}^T)$ , where  $-V_{\mathcal{N}}^T D V_{\mathcal{N}}$  has at least one 0 eigenvalue, correspond to points on the relative boundary of the EDM cone.

$r = 1$ : Symmetric hollow nonnegative matrices rank-one on  $\mathcal{N}(\mathbf{1}^T)$  correspond to extreme directions (1355) of the EDM cone; *id est*, for some nonzero vector  $u$  (§A.3.1.0.7)

$$\left. \begin{array}{l} \text{rank } V_{\mathcal{N}}^T D V_{\mathcal{N}} = 1 \\ D \in \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N} \end{array} \right\} \Leftrightarrow \begin{array}{l} D \in \mathbb{EDM}^N \\ D \text{ is an extreme direction} \end{array} \Leftrightarrow \left\{ \begin{array}{l} -V_{\mathcal{N}}^T D V_{\mathcal{N}} \equiv u u^T \\ D \in \mathbb{S}_h^N \end{array} \right. \quad (1361)$$

**6.5.0.0.1 Proof.** Case  $r = 1$  is easily proved: From the nonnegativity development in §5.8.1, extreme direction (1355), and Schoenberg criterion (1052), we need show only sufficiency; *id est*, prove

$$\left. \begin{array}{l} \text{rank } V_{\mathcal{N}}^T D V_{\mathcal{N}} = 1 \\ D \in \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N} \end{array} \right\} \Rightarrow \begin{array}{l} D \in \mathbb{EDM}^N \\ D \text{ is an extreme direction} \end{array}$$

Any symmetric matrix  $D$  satisfying the rank condition must have the form, for  $z, q \in \mathbb{R}^N$  and nonzero  $z \in \mathcal{N}(\mathbf{1}^T)$ ,

$$D = \pm(\mathbf{1} q^T + q \mathbf{1}^T - 2 z z^T) \quad (1362)$$

because (§5.6.2.1, *confer* §E.7.2.0.2)

$$\mathcal{N}(\mathbf{V}_{\mathcal{N}}(D)) = \{\mathbf{1} q^T + q \mathbf{1}^T \mid q \in \mathbb{R}^N\} \subseteq \mathbb{S}^N \quad (1363)$$

Hollowness demands  $q = \delta(z z^T)$  while nonnegativity demands choice of positive sign in (1362). Matrix  $D$  thus takes the form of an extreme direction (1355) of the EDM cone. ♦

The foregoing proof is not extensible in rank: An EDM with corresponding affine dimension  $r$  has the general form, for  $\{z_i \in \mathcal{N}(\mathbf{1}^T), i=1 \dots r\}$  an independent set,

$$D = \mathbf{1} \delta \left( \sum_{i=1}^r z_i z_i^T \right)^T + \delta \left( \sum_{i=1}^r z_i z_i^T \right) \mathbf{1}^T - 2 \sum_{i=1}^r z_i z_i^T \in \mathbb{EDM}^N \quad (1364)$$

The EDM so defined relies principally on the sum  $\sum z_i z_i^T$  having positive summand coefficients ( $\Leftrightarrow -V_{\mathcal{N}}^T D V_{\mathcal{N}} \succeq 0$ ) 6.8. Then it is easy to find a sum incorporating negative coefficients while meeting rank, nonnegativity, and symmetric hollowness conditions but not positive semidefiniteness on subspace  $\mathcal{R}(V_{\mathcal{N}})$ ; *e.g.*, from page 395,

$$-V \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 5 \\ 1 & 5 & 0 \end{bmatrix} V \frac{1}{2} = z_1 z_1^T - z_2 z_2^T \quad (1365)$$

---

6.8 ( $\Leftarrow$ ) For  $a_i \in \mathbb{R}^{N-1}$ , let  $z_i = V_{\mathcal{N}}^{\dagger T} a_i$ .

**6.5.0.0.2 Example.** *Extreme rays versus rays on the boundary.*

The EDM  $D = \begin{bmatrix} 0 & 1 & 4 \\ 1 & 0 & 1 \\ 4 & 1 & 0 \end{bmatrix}$  is an extreme direction of  $\text{EDM}^3$  where  $u = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$  in (1361).

Because  $-V_N^T D V_N$  has eigenvalues  $\{0, 5\}$ , the ray whose direction is  $D$  also lies on the relative boundary of  $\text{EDM}^3$ .

In exception, EDM  $D = \kappa \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ , for any particular  $\kappa > 0$ , is an extreme direction of  $\text{EDM}^2$  but  $-V_N^T D V_N$  has only one eigenvalue:  $\{\kappa\}$ . Because  $\text{EDM}^2$  is a ray whose relative boundary (§2.6.1.4.1) is the origin, this conventional boundary does not include  $D$  which belongs to the relative interior in this dimension. (§2.7.0.0.1)  $\square$

### 6.5.1 Gram-form correspondence to $\mathbb{S}_+^{N-1}$

With respect to  $\mathbf{D}(G) = \delta(G)\mathbf{1}^T + \mathbf{1}\delta(G)^T - 2G$  (1045) the linear Gram-form EDM operator, results in §5.6.1 provide [3, §2.6]

$$\text{EDM}^N = \mathbf{D}(\mathbf{V}(\text{EDM}^N)) \equiv \mathbf{D}(V_N \mathbb{S}_+^{N-1} V_N^T) \quad (1366)$$

$$V_N \mathbb{S}_+^{N-1} V_N^T \equiv \mathbf{V}(\mathbf{D}(V_N \mathbb{S}_+^{N-1} V_N^T)) = \mathbf{V}(\text{EDM}^N) \triangleq -V \text{EDM}^N V \frac{1}{2} = \mathbb{S}_c^N \cap \mathbb{S}_+^N \quad (1367)$$

a one-to-one correspondence between  $\text{EDM}^N$  and  $\mathbb{S}_+^{N-1}$ .

### 6.5.2 EDM cone by ellipope

Having defined the ellipope parametrized by scalar  $t > 0$

$$\mathcal{E}_t^N = \mathbb{S}_+^N \cap \{\Phi \in \mathbb{S}^N \mid \delta(\Phi) = t\mathbf{1}\} \quad (1242)$$

then following Alfakih [10] we have

$$\text{EDM}^N = \overline{\text{cone}\{\mathbf{1}\mathbf{1}^T - \mathcal{E}_1^N\}} = \overline{\{t(\mathbf{1}\mathbf{1}^T - \mathcal{E}_1^N) \mid t \geq 0\}} \quad (1368)$$

Identification  $\mathcal{E}^N = \mathcal{E}_1^N$  equates the standard ellipope (§5.9.1.0.1, Figure 155) to our parametrized ellipope.

**6.5.2.0.1 Expository.** *Normal cone, tangent cone, ellipope.*

Define  $T_{\mathcal{E}}(\mathbf{1}\mathbf{1}^T)$  to be the *tangent cone* to the ellipope  $\mathcal{E}$  at point  $\mathbf{1}\mathbf{1}^T$ ; *id est*,

$$T_{\mathcal{E}}(\mathbf{1}\mathbf{1}^T) \triangleq \overline{\{t(\mathcal{E} - \mathbf{1}\mathbf{1}^T) \mid t \geq 0\}} \quad (1369)$$

The normal cone  $\mathcal{K}_{\mathcal{E}}^\perp(\mathbf{1}\mathbf{1}^T)$  to the ellipope at  $\mathbf{1}\mathbf{1}^T$  is a closed convex cone defined (§E.10.3.2.1, Figure 203)

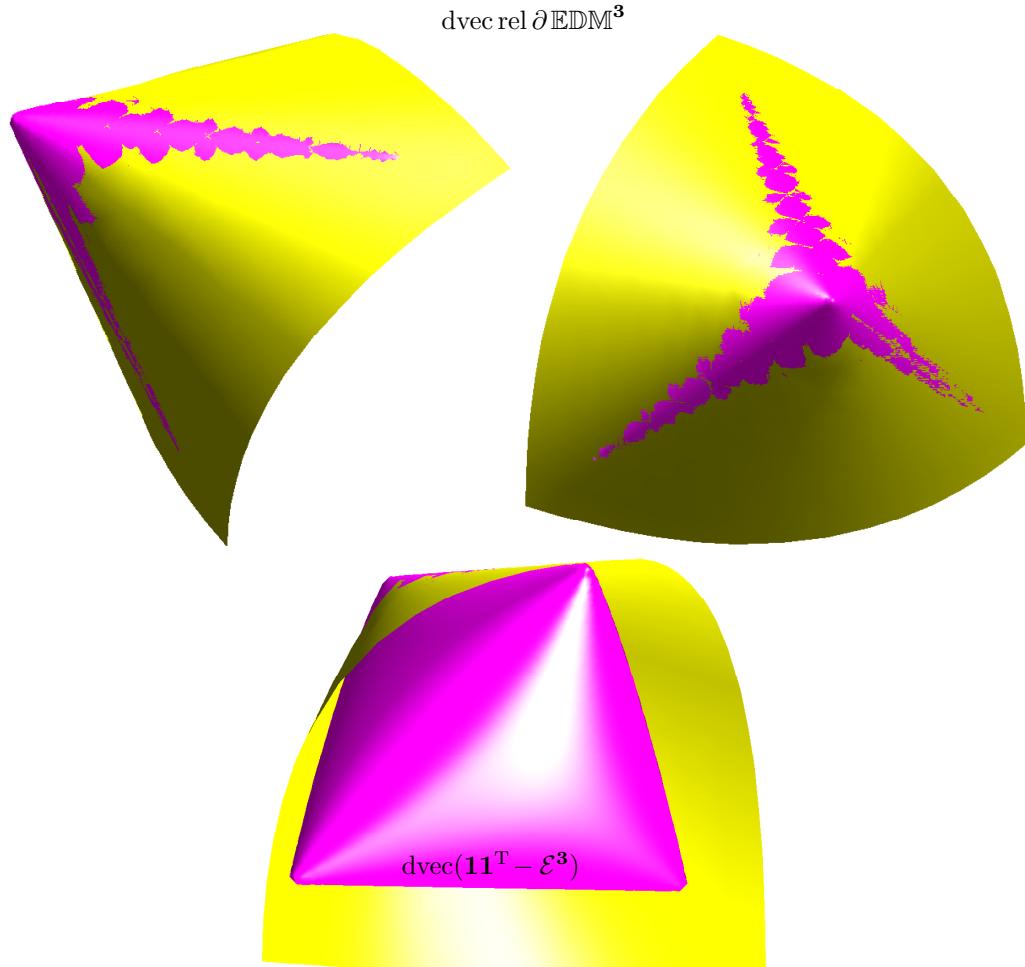
$$\mathcal{K}_{\mathcal{E}}^\perp(\mathbf{1}\mathbf{1}^T) \triangleq \{B \mid \langle B, \Phi - \mathbf{1}\mathbf{1}^T \rangle \leq 0, \Phi \in \mathcal{E}\} \quad (1370)$$

The *polar cone* of any set  $\mathcal{K}$  is the closed convex cone (*confer*(300))

$$\mathcal{K}^\circ \triangleq \{B \mid \langle B, A \rangle \leq 0, \text{ for all } A \in \mathcal{K}\} \quad (1371)$$

The normal cone is well known to be the polar of the tangent cone,

$$\mathcal{K}_{\mathcal{E}}^\perp(\mathbf{1}\mathbf{1}^T) = T_{\mathcal{E}}(\mathbf{1}\mathbf{1}^T)^\circ \quad (1372)$$



$$\mathbb{EDM}^N = \overline{\text{cone}\{\mathbf{1}\mathbf{1}^T - \mathcal{E}^N\}} = \overline{\{t(\mathbf{1}\mathbf{1}^T - \mathcal{E}^N) \mid t \geq 0\}} \quad (1368)$$

Figure 166: Three views of translated negated ellipope  $\mathbf{1}\mathbf{1}^T - \mathcal{E}_1^3$  (*confer* Figure 155) shrouded by truncated EDM cone. Fractal on EDM cone relative boundary is numerical artifact belonging to intersection with ellipope relative boundary. The fractal is trying to convey existence of a neighborhood about the origin where the translated ellipope boundary and EDM cone boundary intersect.

and *vice versa*; [225, §A.5.2.4]

$$\mathcal{K}_{\mathcal{E}}^{\perp}(\mathbf{1}\mathbf{1}^T)^{\circ} = \mathcal{T}_{\mathcal{E}}(\mathbf{1}\mathbf{1}^T) \quad (1373)$$

From Deza & Laurent [126, p.535] we have the EDM cone

$$\mathbb{EDM} = -\mathcal{T}_{\mathcal{E}}(\mathbf{1}\mathbf{1}^T) \quad (1374)$$

The polar EDM cone is also expressible in terms of the ellipope. From (1372) we have

$$\mathbb{EDM}^{\circ} = -\mathcal{K}_{\mathcal{E}}^{\perp}(\mathbf{1}\mathbf{1}^T) \quad (1375)$$

★

In §5.10.1 we proposed the expression for EDM  $D$

$$D = t\mathbf{1}\mathbf{1}^T - \mathfrak{E} \in \mathbb{EDM}^N \quad (1243)$$

where  $t \in \mathbb{R}_+$  and  $\mathfrak{E}$  belongs to the parametrized ellipope  $\mathcal{E}_t^N$ . We further propose, for any particular  $t > 0$

$$\mathbb{EDM}^N = \overline{\text{cone}\{t\mathbf{1}\mathbf{1}^T - \mathcal{E}_t^N\}} \quad (1376)$$

**Proof (pending).** ■

#### 6.5.2.0.2 Exercise. EDM cone from ellipope.

Relationship of the translated negated ellipope with the EDM cone is illustrated in Figure 166. Prove whether it holds that

$$\mathbb{EDM}^N = \overline{\lim_{t \rightarrow \infty} t\mathbf{1}\mathbf{1}^T - \mathcal{E}_t^N} \quad (1377)$$

▼

## 6.6 Vectorization & projection interpretation

In §E.7.2.0.2 we learn:  $-V D V$  can be interpreted as orthogonal projection [7, §2] of vectorized  $-D \in \mathbb{S}_h^N$  on the subspace of geometrically centered symmetric matrices

$$\begin{aligned} \mathbb{S}_c^N &= \{G \in \mathbb{S}^N \mid G\mathbf{1} = \mathbf{0}\} \\ &= \{G \in \mathbb{S}^N \mid \mathcal{N}(G) \supseteq \mathbf{1}\} = \{G \in \mathbb{S}^N \mid \mathcal{R}(G) \subseteq \mathcal{N}(\mathbf{1}^T)\} \\ &= \{VYV \mid Y \in \mathbb{S}^N\} \subset \mathbb{S}^N \\ &\equiv \{V_{\mathcal{N}} A V_{\mathcal{N}}^T \mid A \in \mathbb{S}^{N-1}\} \end{aligned} \quad (1135)$$

because elementary auxiliary matrix  $V$  is an orthogonal projector (§B.4.1). Yet there is another useful projection interpretation:

Revising the fundamental matrix criterion for membership to the EDM cone (1028),<sup>6.9</sup>

$$\left. \begin{array}{l} \langle zz^T, -D \rangle \geq 0 \quad \forall zz^T \mid \mathbf{1}\mathbf{1}^T z z^T = \mathbf{0} \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1378)$$

this is equivalent, of course, to the Schoenberg criterion

$$\left. \begin{array}{l} -V_{\mathcal{N}}^T D V_{\mathcal{N}} \succeq 0 \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1052)$$

<sup>6.9</sup>  $\mathcal{N}(\mathbf{1}\mathbf{1}^T) = \mathcal{N}(\mathbf{1}^T)$  and  $\mathcal{R}(zz^T) = \mathcal{R}(z)$

because  $\mathcal{N}(\mathbf{1}\mathbf{1}^T) = \mathcal{R}(V_N)$ . When  $D \in \mathbb{EDM}^N$ , correspondence (1378) means  $-z^T D z$  is proportional to a nonnegative coefficient of orthogonal projection (§E.6.4.2, Figure 168) of  $-D$  in isometrically isomorphic  $\mathbb{R}^{N(N+1)/2}$  on the range of each and every vectorized (§2.2.2.1) symmetric dyad (§B.1) in the nullspace of  $\mathbf{1}\mathbf{1}^T$ ; *id est*, on each and every member of

$$\begin{aligned}\mathcal{T} &\triangleq \{\text{svec}(zz^T) \mid z \in \mathcal{N}(\mathbf{1}\mathbf{1}^T) = \mathcal{R}(V_N)\} \subset \text{svec } \partial \mathbb{S}_+^N \\ &= \left\{ \text{svec}(V_N v v^T V_N^T) \mid v \in \mathbb{R}^{N-1} \right\}\end{aligned}\quad (1379)$$

whose dimension is

$$\dim \mathcal{T} = N(N-1)/2 \quad (1380)$$

The set of all symmetric dyads  $\{zz^T \mid z \in \mathbb{R}^N\}$  constitute the extreme directions of the positive semidefinite cone (§2.8.1, §2.9)  $\mathbb{S}_+^N$ , hence lie on its boundary. Yet only those dyads in  $\mathcal{R}(V_N)$  are included in the test (1378), thus only a subset  $\mathcal{T}$  of all vectorized extreme directions of  $\mathbb{S}_+^N$  is observed.

In the particularly simple case  $D \in \mathbb{EDM}^2 = \{D \in \mathbb{S}_h^2 \mid d_{12} \geq 0\}$ , for example, only one extreme direction of the PSD cone is involved:

$$zz^T = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \quad (1381)$$

Any nonnegative scaling of vectorized  $zz^T$  belongs to the set  $\mathcal{T}$  illustrated in Figure 167 and Figure 168.

### 6.6.1 Face of PSD cone $\mathbb{S}_+^N$ containing $V$

In any case, set  $\mathcal{T}$  (1379) constitutes the vectorized extreme directions of an  $N(N-1)/2$ -dimensional face of the PSD cone  $\mathbb{S}_+^N$  containing auxiliary matrix  $V$ ; a face isomorphic with  $\mathbb{S}_+^{N-1} = \mathbb{S}_+^{\text{rank } V}$  (§2.9.2.3).

To show this, we must first find the smallest face that contains auxiliary matrix  $V$  and then determine its extreme directions. From (225),

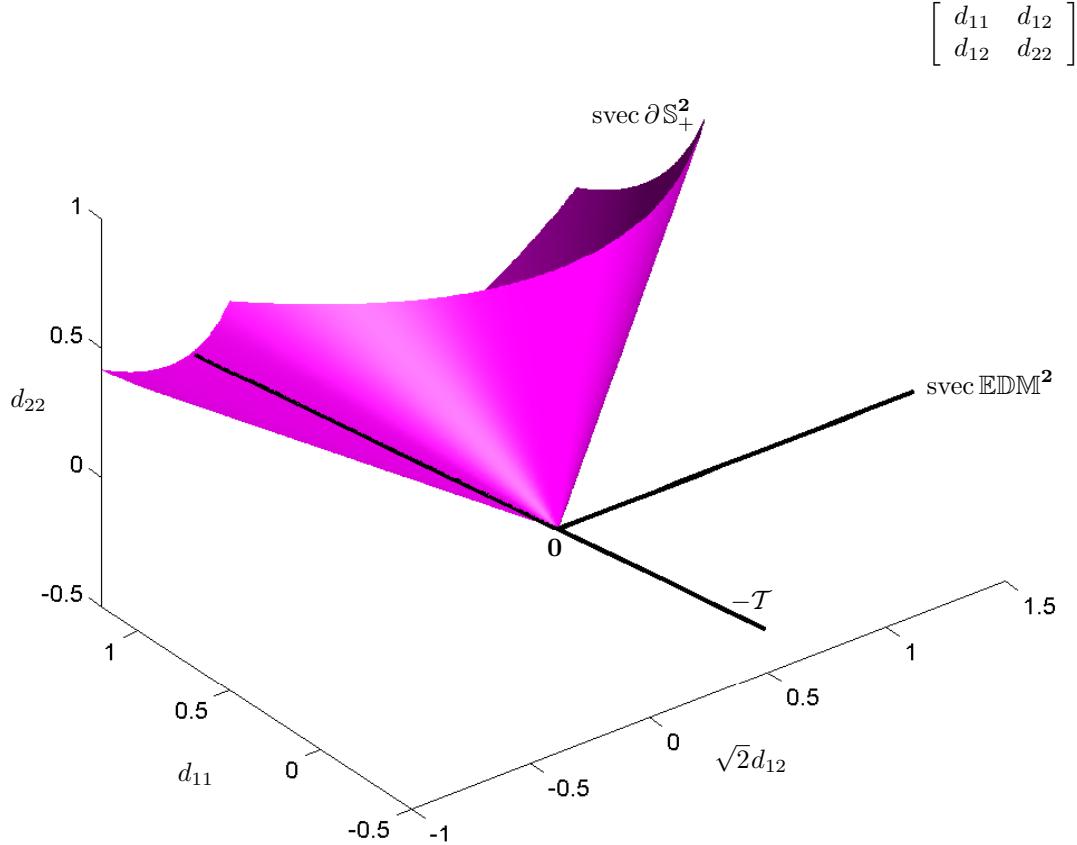
$$\begin{aligned}\mathcal{F}(\mathbb{S}_+^N \ni V) &= \{W \in \mathbb{S}_+^N \mid \mathcal{N}(W) \supseteq \mathcal{N}(V)\} = \{W \in \mathbb{S}_+^N \mid \mathcal{N}(W) \supseteq \mathbf{1}\} \\ &= \{VYV^T \succeq 0 \mid Y \in \mathbb{S}^N\} \equiv \{V_N B V_N^T \mid B \in \mathbb{S}_+^{N-1}\} \\ &\simeq \mathbb{S}_+^{\text{rank } V} = -V_N^T \mathbb{EDM}^N V_N\end{aligned}\quad (1382)$$

where the equivalence  $\equiv$  is from §5.6.1 while isomorphic equality  $\simeq$  with transformed EDM cone is from (1165). Projector  $V$  belongs to  $\mathcal{F}(\mathbb{S}_+^N \ni V)$  because  $V_N V_N^\dagger V_N^{\dagger T} V_N^T = V$  (§B.4.3). Each and every rank-one matrix belonging to this face is therefore of the form:

$$V_N v v^T V_N^T \mid v \in \mathbb{R}^{N-1} \quad (1383)$$

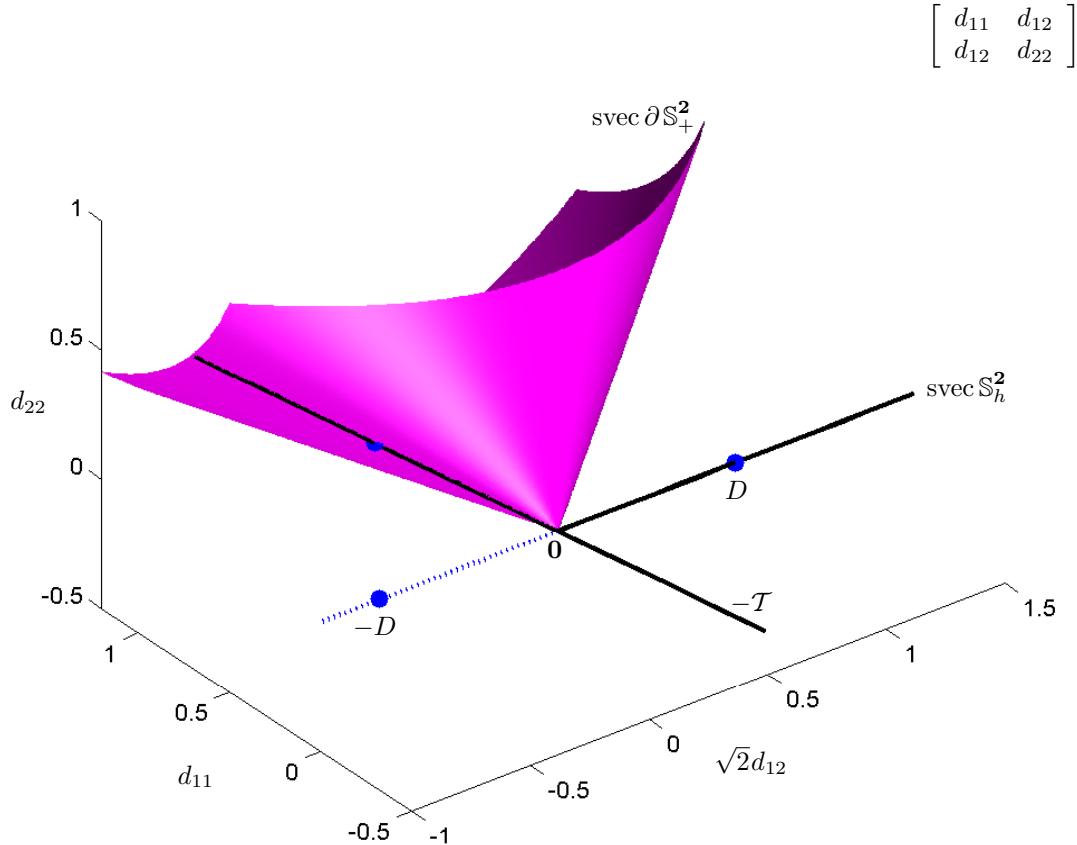
Because  $\mathcal{F}(\mathbb{S}_+^N \ni V)$  is isomorphic with a positive semidefinite cone  $\mathbb{S}_+^{N-1}$ , then  $\mathcal{T}$  constitutes the vectorized extreme directions of  $\mathcal{F}$ , the origin constitutes the extreme points of  $\mathcal{F}$ , and auxiliary matrix  $V$  is some convex combination of those extreme points and directions by the *extremes theorem* (§2.8.1.1.1).





$$\mathcal{T} \triangleq \left\{ \text{svec}(zz^T) \mid z \in \mathcal{N}(\mathbf{1}\mathbf{1}^T) = \kappa \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \kappa \in \mathbb{R} \right\} \subset \text{svec } \partial \mathbb{S}_+^2$$

Figure 167: Truncated boundary of positive semidefinite cone  $\mathbb{S}_+^2$  in isometrically isomorphic  $\mathbb{R}^3$  (via  $\text{svec}$  (57)) is, in this dimension, constituted solely by its extreme directions. Truncated cone of Euclidean distance matrices  $\text{EDM}^2$  drawn in isometrically isomorphic subspace  $\mathbb{R}^3$ . Relative boundary of EDM cone is constituted solely by matrix  $\mathbf{0}$ . Halfline  $\mathcal{T} = \{\kappa^2[1 \ -\sqrt{2} \ 1]^T \mid \kappa \in \mathbb{R}\}$  on PSD cone boundary depicts that lone extreme ray (1381) on which orthogonal projection of  $-D$  must be positive semidefinite if  $D$  is to belong to  $\text{EDM}^2$ . aff cone  $\mathcal{T} = \text{svec } \mathbb{S}_c^2$ . (1386) Dual EDM cone is halfspace in  $\mathbb{R}^3$  whose bounding hyperplane has inward-normal  $\text{svec } \text{EDM}^2$ .



$$P_{\text{svec } zz^T}(\text{svec}(-D)) = \frac{\langle zz^T, -D \rangle}{\langle zz^T, zz^T \rangle} zz^T \text{ is projection of vectorized } -D \text{ on range of vectorized } zz^T.$$

$$D \in \mathbb{EDM}^N \Leftrightarrow \begin{cases} \langle zz^T, -D \rangle \geq 0 \quad \forall zz^T \mid \mathbf{1}\mathbf{1}^T zz^T = \mathbf{0} \\ D \in \mathbb{S}_h^N \end{cases} \quad (1378)$$

Figure 168: Given-matrix  $D$  is assumed to belong to symmetric hollow subspace  $\mathbb{S}_h^2$ ; a line in this dimension. Negative  $D$  is found along  $\mathbb{S}_h^2$ . Set  $\mathcal{T}$  (1379) has only one ray member in this dimension; not orthogonal to  $\mathbb{S}_h^2$ . Orthogonal projection of  $-D$  on  $\mathcal{T}$  (indicated by half dot) has nonnegative projection coefficient. Matrix  $D$  must therefore be an EDM.

In fact the smallest face, that contains auxiliary matrix  $V$ , of the PSD cone  $\mathbb{S}_+^N$  is the intersection with the geometric center subspace (2196) (2197);

$$\begin{aligned}\mathcal{F}(\mathbb{S}_+^N \ni V) &= \text{cone} \left\{ V_N v v^T V_N^T \mid v \in \mathbb{R}^{N-1} \right\} \\ &= \mathbb{S}_c^N \cap \mathbb{S}_+^N \\ &\equiv \{X \succeq 0 \mid \langle X, \mathbf{1}\mathbf{1}^T \rangle = 0\} \quad (1751)\end{aligned}\tag{1384}$$

In isometrically isomorphic  $\mathbb{R}^{N(N+1)/2}$

$$\text{svec } \mathcal{F}(\mathbb{S}_+^N \ni V) = \text{cone } \mathcal{T} \tag{1385}$$

related to  $\mathbb{S}_c^N$  by

$$\text{aff cone } \mathcal{T} = \text{svec } \mathbb{S}_c^N \tag{1386}$$

### 6.6.2 EDM criteria in $\mathbf{1}\mathbf{1}^T$

(confer §6.4, (1059)) Laurent specifies an ellipope trajectory condition for EDM cone membership: [264, §2.3]

$$D \in \mathbb{EDM}^N \Leftrightarrow [1 - e^{-\alpha d_{ij}}] \in \mathbb{EDM}^N \quad \forall \alpha > 0 \quad (1237a)$$

From the parametrized ellipope  $\mathcal{E}_t^N$  in §6.5.2 and §5.10.1 we propose

$$D \in \mathbb{EDM}^N \Leftrightarrow \exists \begin{cases} t \in \mathbb{R}_+ \\ \mathfrak{E} \in \mathcal{E}_t^N \end{cases} \ni D = t\mathbf{1}\mathbf{1}^T - \mathfrak{E} \tag{1387}$$

Chabriac & Crouzeix [81, §4] prove a different criterion they attribute to Finsler, 1937 [163]. We apply it to EDMs: for  $D \in \mathbb{S}_h^N$  (1185)

$$\begin{aligned}-V_N^T D V_N \succ 0 &\Leftrightarrow \exists \kappa > 0 \ni -D + \kappa \mathbf{1}\mathbf{1}^T \succ 0 \\ D \in \mathbb{EDM}^N \text{ with corresponding affine dimension } r = N-1\end{aligned}\tag{1388}$$

This *Finsler criterion* has geometric interpretation in terms of the vectorization & projection already discussed in connection with (1378). With reference to Figure 167, the offset  $\mathbf{1}\mathbf{1}^T$  is simply a direction orthogonal to  $\mathcal{T}$  in isomorphic  $\mathbb{R}^3$ . Intuitively, translation of  $-D$  in direction  $\mathbf{1}\mathbf{1}^T$  is like orthogonal projection on  $\mathcal{T}$  insofar as similar information can be obtained.

When the Finsler criterion (1388) is applied despite lower affine dimension, the constant  $\kappa$  can go to infinity making the test  $-D + \kappa \mathbf{1}\mathbf{1}^T \succeq 0$  impractical for numerical computation. Chabriac & Crouzeix invent a criterion for the semidefinite case, but is no more practical: for  $D \in \mathbb{S}_h^N$

$$\begin{aligned}D \in \mathbb{EDM}^N &\Leftrightarrow \\ \exists \kappa_p > 0 \ni \forall \kappa \geq \kappa_p, -D - \kappa \mathbf{1}\mathbf{1}^T [\text{sic}] \text{ has exactly one negative eigenvalue}\end{aligned}\tag{1389}$$

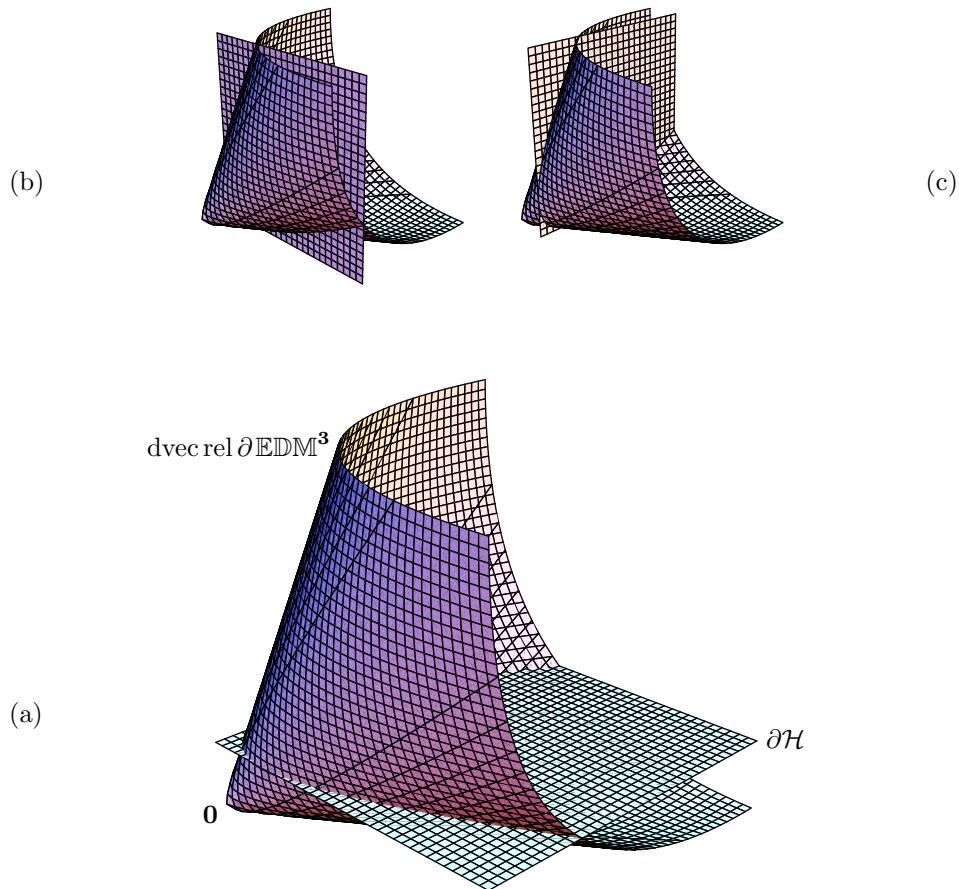


Figure 169: (a) In isometrically isomorphic subspace  $\mathbb{R}^3$ , intersection of  $\text{EDM}^3$  with hyperplane  $\partial H$  representing one fixed symmetric entry  $d_{23} = \kappa$  (both drawn truncated, rounded vertex is artifact of plot). EDMs in this dimension corresponding to affine dimension 1 comprise relative boundary of EDM cone (§6.5). Since intersection illustrated includes a nontrivial subset of cone's relative boundary, then it is apparent there exist infinitely many EDM completions corresponding to affine dimension 1. In this dimension it is impossible to represent a unique nonzero completion corresponding to affine dimension 1, for example, using a single hyperplane because any hyperplane supporting relative boundary at a particular point  $\Gamma$  contains an entire ray  $\{\zeta \Gamma \mid \zeta \geq 0\}$  belonging to  $\text{rel } \partial \text{EDM}^3$  by Lemma 2.8.0.0.1. (b)  $d_{13} = \kappa$ . (c)  $d_{12} = \kappa$ .

## 6.7 A geometry of completion

*It is not known how to proceed if one wishes to restrict the dimension of the Euclidean space in which the configuration of points may be constructed.*

— Michael W. Trosset, 2000 [393, §1]

Given an incomplete noiseless EDM, intriguing is the question of whether a list in  $X \in \mathbb{R}^{n \times N}$  (77) may be reconstructed and under what circumstances reconstruction is unique. [3] [5] [6] [7] [9] [18] [73] [232] [244] [263] [264] [265]

If one or more entries of a particular EDM are fixed, then geometric interpretation of the feasible set of completions is the intersection of the EDM cone  $\text{EDM}^N$  in isomorphic subspace  $\mathbb{R}^{N(N-1)/2}$  with as many hyperplanes as there are fixed symmetric entries. 6.10 Assuming a nonempty intersection, then the number of completions is generally infinite, and those corresponding to particular affine dimension  $r < N - 1$  belong to some generally nonconvex subset of that intersection (*confer* §2.9.2.9.2) that can be as small as a point.

### 6.7.0.0.1 Example. Maximum variance unfolding.

[435]

A process minimizing affine dimension (§2.1.5) of certain kinds of Euclidean manifold by topological transformation can be posed as a completion problem (*confer* §E.10.2.1.2). Weinberger & Saul, who originated the technique, specify an applicable manifold in three dimensions by analogy to an ordinary sheet of paper (*confer* §2.1.6); imagine, we find it deformed from flatness in some way introducing neither holes, tears, or selfintersections. [435, §2.2] The physical process is intuitively described as *unfurling*, *unfolding*, *diffusing*, *decompacting*, or *unraveling*. In particular instances, the process is a sort of flattening by stretching until taut (but not by crushing); *e.g.*, unfurling a three-dimensional Euclidean body resembling a billowy national flag reduces that manifold's affine dimension to  $r=2$ .

Data input, to the proposed process, originates from distances between relatively dense neighboring samples of a given manifold. Figure 170 realizes a densely sampled neighborhood; called, *neighborhood graph*. Essentially, the algorithmic process preserves local isometry between *nearest neighbors* allowing distant neighbors to excursion expansively by “maximizing variance” (Figure 7). A common number, of nearest neighbors to each sample, is a data-dependent algorithmic parameter whose minimum value connects the graph. A dimensionless *EDM subgraph*, between each sample and its nearest neighbors, is completed from available data then included as input. One such EDM subgraph completion is drawn superimposed upon the neighborhood graph in Figure 170. 6.11 The consequent dimensionless EDM graph, comprising all subgraphs, is generally incomplete because neighbor number is relatively small; incomplete though it is a superset of the neighborhood graph. Remaining distances (those not graphed at all) are squared then made variables within the algorithm. It is this variability that admits unfurling.

To demonstrate, consider untying the *trefoil knot* drawn in Figure 171a. A corresponding EDM  $D = [d_{ij}, i, j = 1 \dots N]$  employing only two nearest neighbors is banded having the incomplete form:

6.10 Depicted in Figure 169a is an intersection of the EDM cone  $\text{EDM}^3$  with a single hyperplane representing the set of all EDMs having one fixed symmetric entry. This representation is practical because it is easily combined with or replaced by other convex constraints; *e.g.*, slab inequalities in (816) that admit bounding of noise processes.

6.11 Local reconstruction of point position, from the EDM submatrix corresponding to a complete dimensionless EDM subgraph, is unique to within an isometry (§5.6, §5.12).

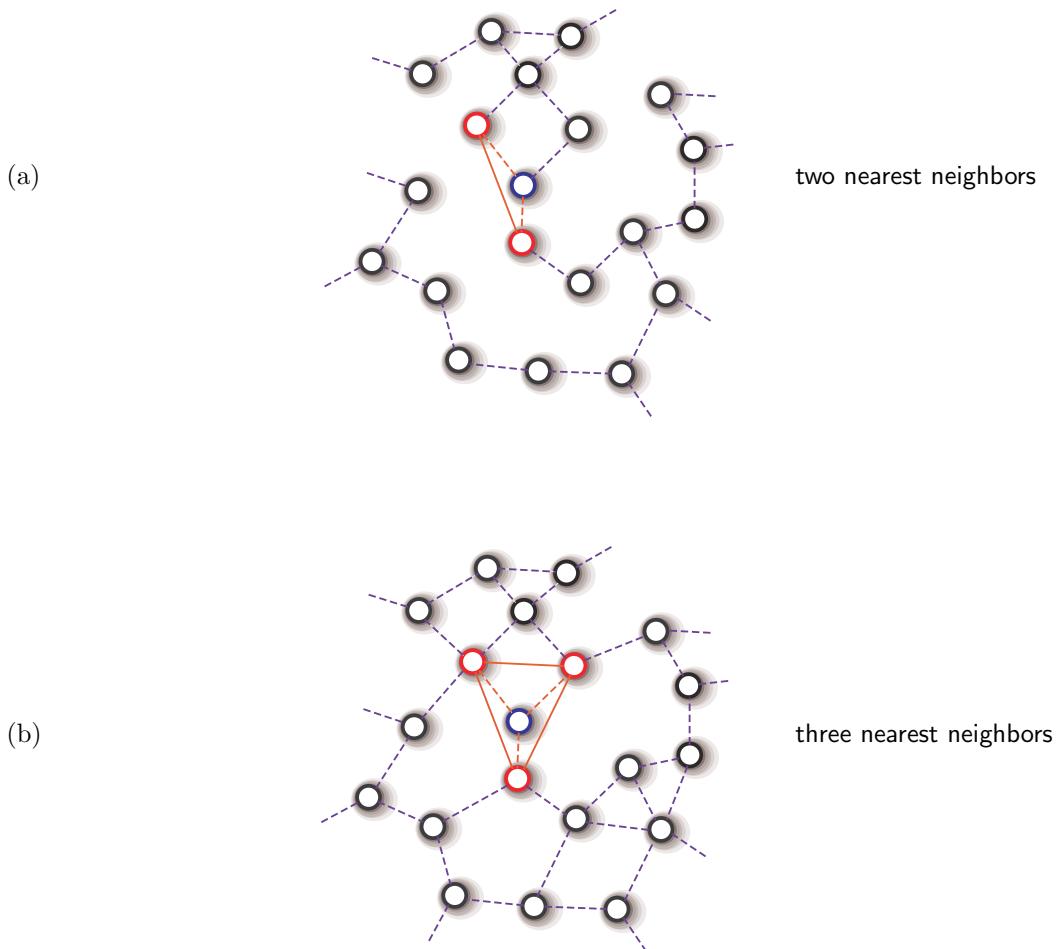


Figure 170: One dimensionless EDM subgraph completion (solid) superimposed on (but not obscuring) neighborhood graph (dashed). Local view of a few dense samples  $\circ$  from relative interior of some arbitrary Euclidean manifold whose affine dimension appears two-dimensional in this neighborhood. All line segments measure absolute distance. Dashed line segments help visually locate nearest neighbors; suggesting, best number of nearest neighbors can be greater than value of embedding dimension after topological transformation (*confer* [239, §2]). Solid line segments represent completion of EDM subgraph from available distance data for an arbitrarily chosen sample and its nearest neighbors. Each distance from EDM subgraph becomes distance-square in corresponding EDM submatrix.

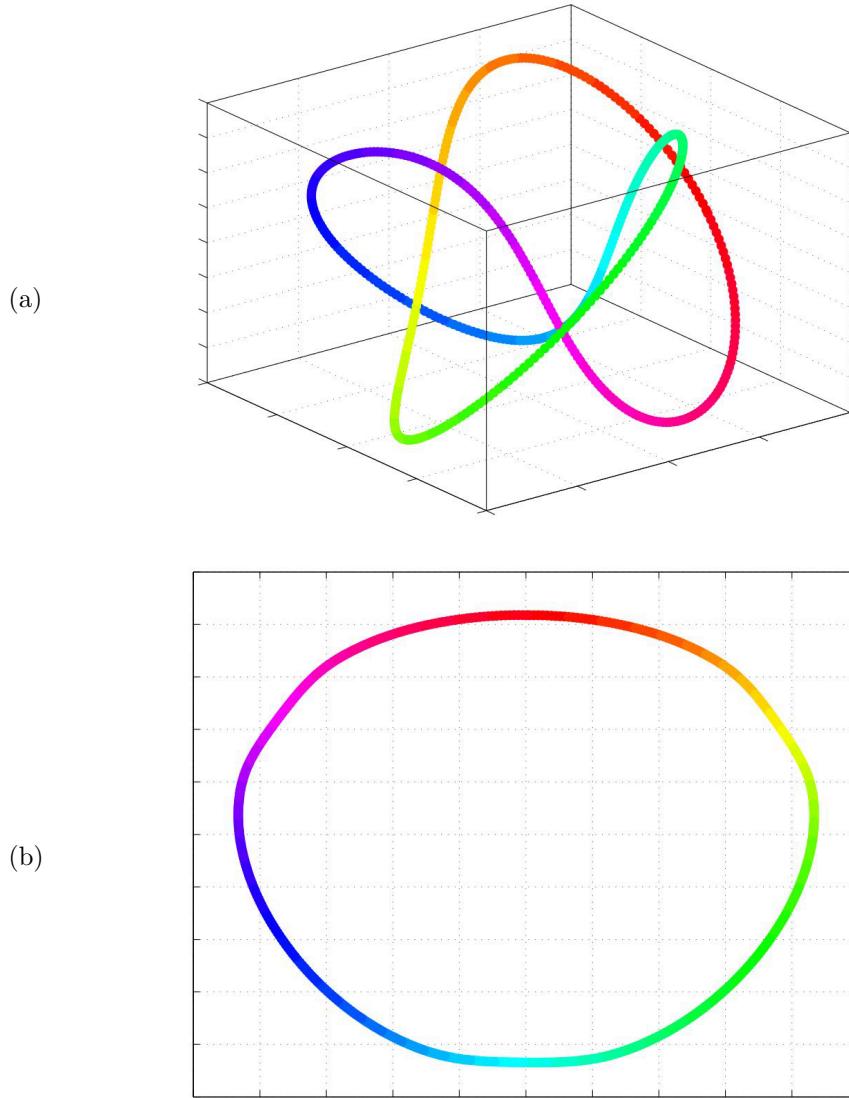


Figure 171: (a) Trefoil knot in  $\mathbb{R}^3$  from Weinberger & Saul [435]. (b) Topological transformation algorithm employing four nearest neighbors and  $N=539$  samples reduces affine dimension of knot to  $r=2$ . Choosing instead two nearest neighbors would make this embedding more circular.

$$D = \begin{bmatrix} 0 & \check{d}_{12} & \check{d}_{13} & ? & \cdots & ? & \check{d}_{1,N-1} & \check{d}_{1N} \\ \check{d}_{12} & 0 & \check{d}_{23} & \check{d}_{24} & \ddots & ? & ? & \check{d}_{2N} \\ \check{d}_{13} & \check{d}_{23} & 0 & \check{d}_{34} & \ddots & ? & ? & ? \\ ? & \check{d}_{24} & \check{d}_{34} & 0 & \ddots & \ddots & ? & ? \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & ? \\ ? & ? & ? & \ddots & \ddots & 0 & \check{d}_{N-2,N-1} & \check{d}_{N-2,N} \\ \check{d}_{1,N-1} & ? & ? & ? & \ddots & \check{d}_{N-2,N-1} & 0 & \check{d}_{N-1,N} \\ \check{d}_{1N} & \check{d}_{2N} & ? & ? & ? & \check{d}_{N-2,N} & \check{d}_{N-1,N} & 0 \end{bmatrix} \quad (1390)$$

where  $\check{d}_{ij}$  denotes a given fixed distance-square. The unfurling algorithm can be expressed as an optimization problem; constrained total distance-square maximization:

$$\begin{aligned} & \underset{D}{\text{maximize}} \quad \langle -V, D \rangle \\ & \text{subject to} \quad \langle D, e_i e_j^T + e_j e_i^T \rangle \frac{1}{2} = \check{d}_{ij} \quad \forall (i, j) \in \mathcal{I} \\ & \qquad \text{rank}(VDV) = 2 \\ & \qquad D \in \mathbb{EDM}^N \end{aligned} \quad (1391)$$

where  $e_i \in \mathbb{R}^N$  is the  $i^{\text{th}}$  member of the standard basis, where set  $\mathcal{I}$  indexes the given distance-square data like that in (1390), where  $V \in \mathbb{R}^{N \times N}$  is the geometric centering matrix (§B.4.1), and where

$$\langle -V, D \rangle = \text{tr}(-VDV) = 2 \text{tr} G = \frac{1}{N} \sum_{i,j} d_{ij} \quad (1057)$$

where  $G$  is the Gram matrix producing  $D$  assuming  $G\mathbf{1} = \mathbf{0}$ .

If the (rank) constraint on affine dimension is ignored, then problem (1391) becomes convex, a corresponding solution  $D^*$  can be found, and a nearest rank-2 solution is then had by ordered eigenvalue decomposition of  $-VD^*V$  followed by *spectral projection* (§7.1.3) on  $\begin{bmatrix} \mathbb{R}_+^2 \\ \mathbf{0} \end{bmatrix} \subset \mathbb{R}^N$ . This two-step process is necessarily suboptimal. Yet because the decomposition for the trefoil knot reveals only two dominant eigenvalues, the spectral projection is nearly benign. Such a reconstruction of point position (§5.12) utilizing four nearest neighbors is drawn in Figure 171b; a low-dimensional embedding of the trefoil knot.

This problem (1391) can, of course, be written equivalently in terms of Gram matrix  $G$ , facilitated by (1063); *videlicet*, for  $\Phi_{ij}$  as in (1031)

$$\begin{aligned} & \underset{G \in \mathbb{S}_c^N}{\text{maximize}} \quad \langle I, G \rangle \\ & \text{subject to} \quad \langle G, \Phi_{ij} \rangle = \check{d}_{ij} \quad \forall (i, j) \in \mathcal{I} \\ & \qquad \text{rank } G = 2 \\ & \qquad G \succeq 0 \end{aligned} \quad (1392)$$

a trace maximization. The advantage to converting EDM to Gram is: Gram matrix  $G$  is a bridge between point list  $X$  and EDM  $D$ ; constraints on any or all of these three variables may now be introduced. (Example 5.4.2.2.8) Confinement to geometric center subspace  $\mathbb{S}_c^N$  (implicit constraint  $G\mathbf{1} = \mathbf{0}$ ) keeps  $G$  independent of  $\mathbb{S}_c^{N \perp}$  its translation-invariant subspace (§5.5.1.1, Figure 173) so as not to become numerically unbounded.

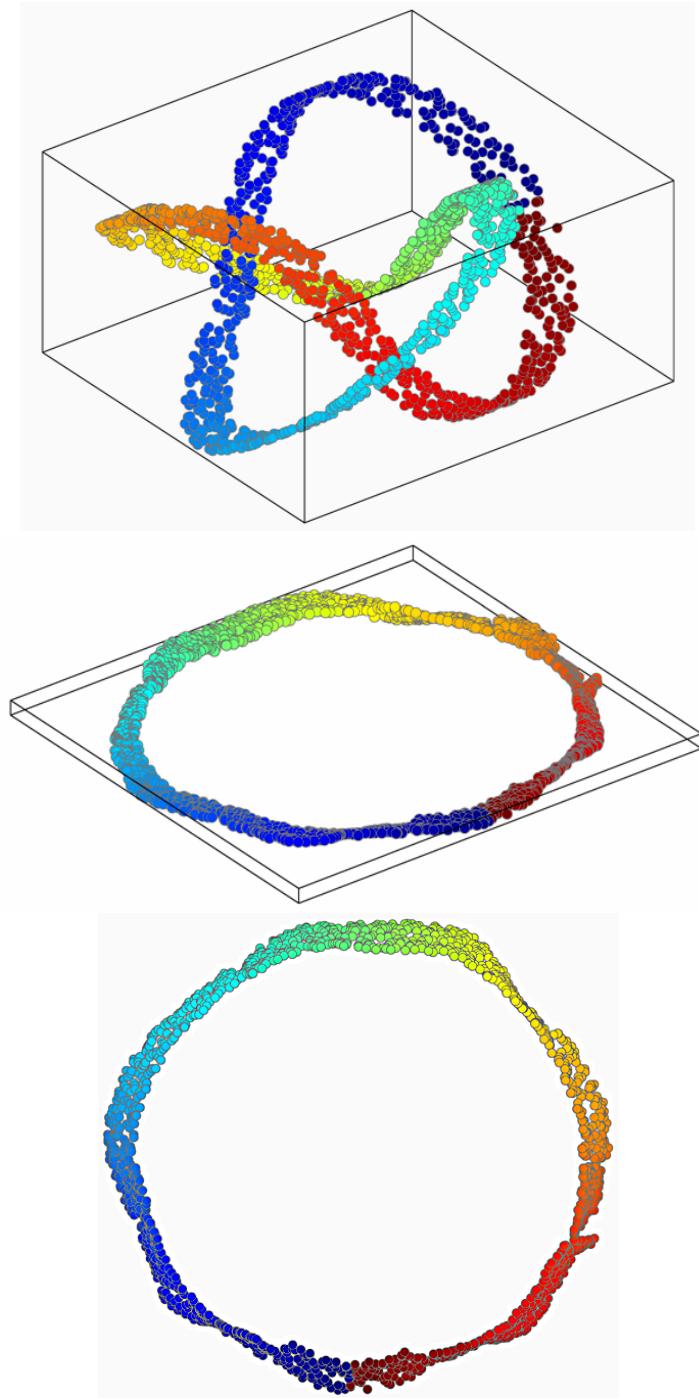


Figure 172: *Trefoil ribbon* (Kilian Weinberger). Same topological transformation algorithm as in Figure 171b with five nearest neighbors and  $N = 1617$  samples.

A problem dual to *maximum variance unfolding problem* (1392) (less the Gram rank constraint) has been called the *fastest mixing Markov process*. That dual has simple interpretations in graph and circuit theory and in mechanical and thermal systems, explored in [377], and has direct application to quick calculation of *PageRank* by search engines [259, §4]. Optimal Gram rank turns out to be tightly bounded above by minimum multiplicity of the second smallest eigenvalue of a dual optimal variable.  $\square$

## 6.8 Dual EDM cone

### 6.8.1 Ambient $\mathbb{S}^N$

We consider finding the ordinary dual EDM cone in ambient space  $\mathbb{S}^N$  where  $\mathbb{EDM}^N$  is pointed, closed, convex, but not full-dimensional. The set of all EDMs in  $\mathbb{S}^N$  is a closed convex cone because it is the intersection of halfspaces about the origin in vectorized variable  $D$  (§2.4.1.1.1, §2.7.2):

$$\mathbb{EDM}^N = \bigcap_{\substack{z \in \mathcal{N}(\mathbf{1}^T) \\ i=1 \dots N}} \left\{ D \in \mathbb{S}^N \mid \langle e_i e_i^T, D \rangle \geq 0, \langle e_i e_i^T, D \rangle \leq 0, \langle z z^T, -D \rangle \geq 0 \right\} \quad (1393)$$

By definition (300), dual cone  $\mathcal{K}^*$  comprises each and every vector inward-normal to a hyperplane supporting convex cone  $\mathcal{K}$ . The dual EDM cone in the ambient space of symmetric matrices is therefore expressible as the aggregate of every conic combination of inward-normals from (1393):

$$\begin{aligned} \mathbb{EDM}^{N^*} &= \text{cone}\{e_i e_i^T, -e_j e_j^T \mid i, j = 1 \dots N\} - \text{cone}\{z z^T \mid \mathbf{1} \mathbf{1}^T z z^T = \mathbf{0}\} \\ &= \left\{ \sum_{i=1}^N \zeta_i e_i e_i^T - \sum_{j=1}^N \xi_j e_j e_j^T \mid \zeta_i, \xi_j \geq 0 \right\} - \text{cone}\{z z^T \mid \mathbf{1} \mathbf{1}^T z z^T = \mathbf{0}\} \\ &= \{\delta(u) \mid u \in \mathbb{R}^N\} - \text{cone}\{V_N v v^T V_N^T \mid v \in \mathbb{R}^{N-1}, (\|v\|=1)\} \subset \mathbb{S}^N \\ &= \{\delta^2(Y) - V_N \Psi V_N^T \mid Y \in \mathbb{S}^N, \Psi \in \mathbb{S}_+^{N-1}\} \end{aligned} \quad (1394)$$

The EDM cone is not selfdual in ambient  $\mathbb{S}^N$  because its affine hull belongs to a proper subspace

$$\text{aff } \mathbb{EDM}^N = \mathbb{S}_h^N \quad (1395)$$

The ordinary dual EDM cone cannot, therefore, be pointed. (§2.13.1.2)

When  $N=1$ , the EDM cone is the point at the origin in  $\mathbb{R}$ . Auxiliary matrix  $V_N$  is empty  $[\emptyset]$ , and dual cone  $\mathbb{EDM}^*$  is the real line.

When  $N=2$ , the EDM cone is a nonnegative real line in isometrically isomorphic  $\mathbb{R}^3$ ; there  $\mathbb{S}_h^2$  is a real line containing the EDM cone. Dual cone  $\mathbb{EDM}^{2^*}$  is the particular halfspace in  $\mathbb{R}^3$  whose boundary has inward-normal  $\mathbb{EDM}^2$ . Diagonal matrices  $\{\delta(u)\}$  in (1394) are represented by a hyperplane through the origin  $\{\underline{d} \mid [0 \ 1 \ 0] \underline{d} = 0\}$  while the term  $\text{cone}\{V_N v v^T V_N^T\}$  is represented by the halffline  $\mathcal{T}$  in Figure 167 belonging to the positive semidefinite cone boundary. The dual EDM cone is formed by translating the hyperplane along the negative semidefinite halffline  $-\mathcal{T}$ ; the union of each and every translation. (*confer* §2.10.2.0.1)

When cardinality  $N$  exceeds 2, the dual EDM cone can no longer be polyhedral simply because the EDM cone cannot. (§2.13.1.2)

### 6.8.1.1 EDM cone and its dual in ambient $\mathbb{S}^N$

Consider the two convex cones

$$\begin{aligned}\mathcal{K}_1 &\triangleq \mathbb{S}_h^N \\ \mathcal{K}_2 &\triangleq \bigcap_{y \in \mathcal{N}(\mathbf{1}^T)} \left\{ A \in \mathbb{S}^N \mid \langle yy^T, -A \rangle \geq 0 \right\} \\ &= \left\{ A \in \mathbb{S}^N \mid -z^T V A V z \geq 0 \quad \forall z z^T (\succeq 0) \right\} \\ &= \left\{ A \in \mathbb{S}^N \mid -V A V \succeq 0 \right\}\end{aligned}\tag{1396}$$

so

$$\mathcal{K}_1 \cap \mathcal{K}_2 = \text{EDM}^N\tag{1397}$$

Dual cone  $\mathcal{K}_1^* = \mathbb{S}_h^{N \perp} \subseteq \mathbb{S}^N$  (73) is the subspace of diagonal matrices. From (1394) via (319),

$$\mathcal{K}_2^* = -\text{cone}\left\{ V_N v v^T V_N^T \mid v \in \mathbb{R}^{N-1} \right\} \subset \mathbb{S}^N\tag{1398}$$

Gaffke & Mathar [168, §5.3] observe that projection on  $\mathcal{K}_1$  and  $\mathcal{K}_2$  have simple closed forms: Projection on subspace  $\mathcal{K}_1$  is easily performed by symmetrization and zeroing the main diagonal or *vice versa*, while projection of  $H \in \mathbb{S}^N$  on  $\mathcal{K}_2$  is <sup>6.12</sup>

$$P_{\mathcal{K}_2} H = H - P_{\mathbb{S}_+^N}(V H V)\tag{1399}$$

**Proof.** First, we observe membership of  $H - P_{\mathbb{S}_+^N}(V H V)$  to  $\mathcal{K}_2$  because

$$P_{\mathbb{S}_+^N}(V H V) - H = \left( P_{\mathbb{S}_+^N}(V H V) - V H V \right) + (V H V - H)\tag{1400}$$

The term  $P_{\mathbb{S}_+^N}(V H V) - V H V$  necessarily belongs to the (dual) positive semidefinite cone by Theorem E.9.2.0.1.  $V^2 = V$ , hence

$$-V \left( H - P_{\mathbb{S}_+^N}(V H V) \right) V \succeq 0\tag{1401}$$

by Corollary A.3.1.0.5.

Next, we require

$$\langle P_{\mathcal{K}_2} H - H, P_{\mathcal{K}_2} H \rangle = 0\tag{1402}$$

Expanding,

$$\langle -P_{\mathbb{S}_+^N}(V H V), H - P_{\mathbb{S}_+^N}(V H V) \rangle = 0\tag{1403}$$

$$\langle P_{\mathbb{S}_+^N}(V H V), (P_{\mathbb{S}_+^N}(V H V) - V H V) + (V H V - H) \rangle = 0\tag{1404}$$

$$\langle P_{\mathbb{S}_+^N}(V H V), (V H V - H) \rangle = 0\tag{1405}$$

Product  $V H V$  belongs to the geometric center subspace; (§E.7.2.0.2)

$$V H V \in \mathbb{S}_c^N = \{Y \in \mathbb{S}^N \mid \mathcal{N}(Y) \supseteq \mathbf{1}\}\tag{1406}$$

Diagonalize  $V H V \triangleq Q \Lambda Q^T$  (§A.5) whose nullspace is spanned by the eigenvectors corresponding to 0 eigenvalues by Theorem A.7.3.0.1. Projection of  $V H V$  on the PSD cone (§7.1) simply zeroes negative eigenvalues in diagonal matrix  $\Lambda$ . Then

$$\mathcal{N}(P_{\mathbb{S}_+^N}(V H V)) \supseteq \mathcal{N}(V H V) (\supseteq \mathcal{N}(V))\tag{1407}$$

---

<sup>6.12</sup>  $P_{\mathbb{S}_+^N}(V H V) = \mathbf{0}$  for  $H \in \text{EDM}^N$ .

from which it follows:

$$P_{\mathbb{S}_+^N}(V H V) \in \mathbb{S}_c^N \quad (1408)$$

so  $P_{\mathbb{S}_+^N}(V H V) \perp (V H V - H)$  because  $V H V - H \in \mathbb{S}_c^{N \perp}$ .

Finally, we must have  $P_{\mathcal{K}_2} H - H = -P_{\mathbb{S}_+^N}(V H V) \in \mathcal{K}_2^*$ . Dual cone  $\mathcal{K}_2^* = -\mathcal{F}(\mathbb{S}_+^N \ni V)$  is the negative of the positive semidefinite cone's smallest face that contains auxiliary matrix  $V$ . (§6.6.1) Thus  $P_{\mathbb{S}_+^N}(V H V) \in \mathcal{F}(\mathbb{S}_+^N \ni V) \Leftrightarrow \mathcal{N}(P_{\mathbb{S}_+^N}(V H V)) \supseteq \mathcal{N}(V)$  (§2.9.2.3) which was already established in (1407). ♦

From results in §E.7.2.0.2 we know: matrix product  $V H V = P_{\mathbb{S}_c^N} H$  is the orthogonal projection of  $H \in \mathbb{S}^N$  on the geometric center subspace  $\mathbb{S}_c^N$ . Thus the projection product

$$P_{\mathcal{K}_2} H = H - P_{\mathbb{S}_+^N} P_{\mathbb{S}_c^N} H \quad (1409)$$

**6.8.1.1.1 Lemma.** *Projection on PSD cone ∩ geometric center subspace.*

$$P_{\mathbb{S}_+^N \cap \mathbb{S}_c^N} = P_{\mathbb{S}_+^N} P_{\mathbb{S}_c^N} \quad (1410)$$

◊

**Proof.** For each and every  $H \in \mathbb{S}^N$ , projection of  $P_{\mathbb{S}_c^N} H$  on the positive semidefinite cone remains in the geometric center subspace

$$\begin{aligned} \mathbb{S}_c^N &= \{G \in \mathbb{S}^N \mid G\mathbf{1} = \mathbf{0}\} \\ &= \{G \in \mathbb{S}^N \mid \mathcal{N}(G) \supseteq \mathbf{1}\} = \{G \in \mathbb{S}^N \mid \mathcal{R}(G) \subseteq \mathcal{N}(\mathbf{1}^T)\} \quad (1135) \\ &= \{VYV \mid Y \in \mathbb{S}^N\} \subset \mathbb{S}^N \end{aligned}$$

That is because: eigenvectors of  $P_{\mathbb{S}_c^N} H$ , corresponding to 0 eigenvalues, span its nullspace  $\mathcal{N}(P_{\mathbb{S}_c^N} H)$ . (§A.7.3.0.1) To project  $P_{\mathbb{S}_c^N} H$  on the positive semidefinite cone, its negative eigenvalues are zeroed. (§7.1.2) The nullspace is thereby expanded while eigenvectors originally spanning  $\mathcal{N}(P_{\mathbb{S}_c^N} H)$  remain intact. Because the geometric center subspace is invariant to projection on the PSD cone, then the rule for order of projection on a convex set in a subspace governs (§E.9.5) and statement (1410) follows. ♦

From the lemma it follows

$$\{P_{\mathbb{S}_+^N} P_{\mathbb{S}_c^N} H \mid H \in \mathbb{S}^N\} = \{P_{\mathbb{S}_+^N \cap \mathbb{S}_c^N} H \mid H \in \mathbb{S}^N\} \quad (1411)$$

Then from (2225)

$$-(\mathbb{S}_c^N \cap \mathbb{S}_+^N)^* = \{H - P_{\mathbb{S}_+^N} P_{\mathbb{S}_c^N} H \mid H \in \mathbb{S}^N\} \quad (1412)$$

From (319) we get closure of a vector sum

$$\mathcal{K}_2 = -(\mathbb{S}_c^N \cap \mathbb{S}_+^N)^* = \mathbb{S}_c^{N \perp} - \mathbb{S}_+^N \quad (1413)$$

therefore the equality [112]

$$\text{EDM}^N = \mathcal{K}_1 \cap \mathcal{K}_2 = \mathbb{S}_h^N \cap (\mathbb{S}_c^{N \perp} - \mathbb{S}_+^N) \quad (1414)$$

whose veracity is intuitively evident, in hindsight, [97, p.109] from the most fundamental EDM definition (1033).<sup>6.13</sup> A realization of this construction in low dimension is illustrated in Figure 173 and Figure 174.

<sup>6.13</sup> Formula (1414) is not a matrix criterion for membership to the EDM cone, it is not an EDM definition, and it is not an equivalence between EDM operators or an isomorphism. Rather, it is a recipe for constructing the EDM cone whole from large Euclidean bodies: the positive semidefinite cone, orthogonal complement of the geometric center subspace, and symmetric hollow subspace.

$$\text{EDM}^2 = \mathbb{S}_h^2 \cap \left( \mathbb{S}_c^{2\perp} - \mathbb{S}_+^2 \right)$$

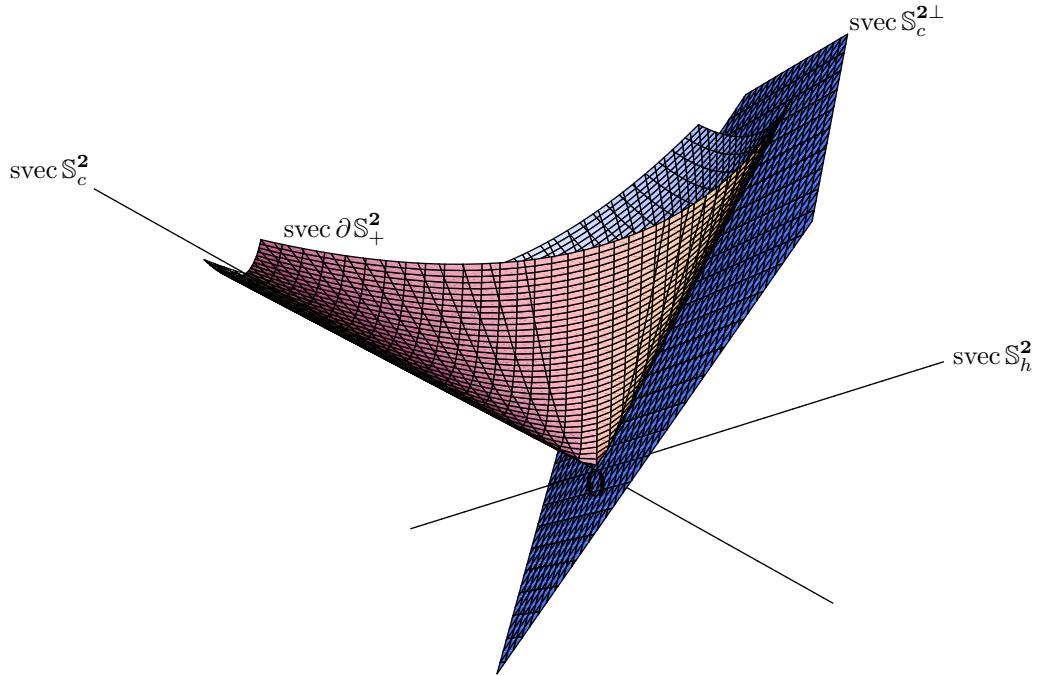


Figure 173: A plane in isometrically isomorphic  $\mathbb{R}^3$ , orthogonal complement  $\mathbb{S}_c^{2\perp}$  (2198) ([§B.2](#)) of geometric center subspace (tiled fragment drawn) apparently supports PSD cone (rounded vertex is plot artifact). Line  $\text{svec } \mathbb{S}_c^2 = \text{aff cone } T$  (1386), intersecting  $\text{svec } \partial \mathbb{S}_+^2$  and drawn in Figure 167, runs along PSD cone boundary. (*confer* Figure 154)

$$\text{EDM}^2 = \mathbb{S}_h^2 \cap \left( \mathbb{S}_c^{2\perp} - \mathbb{S}_+^2 \right)$$

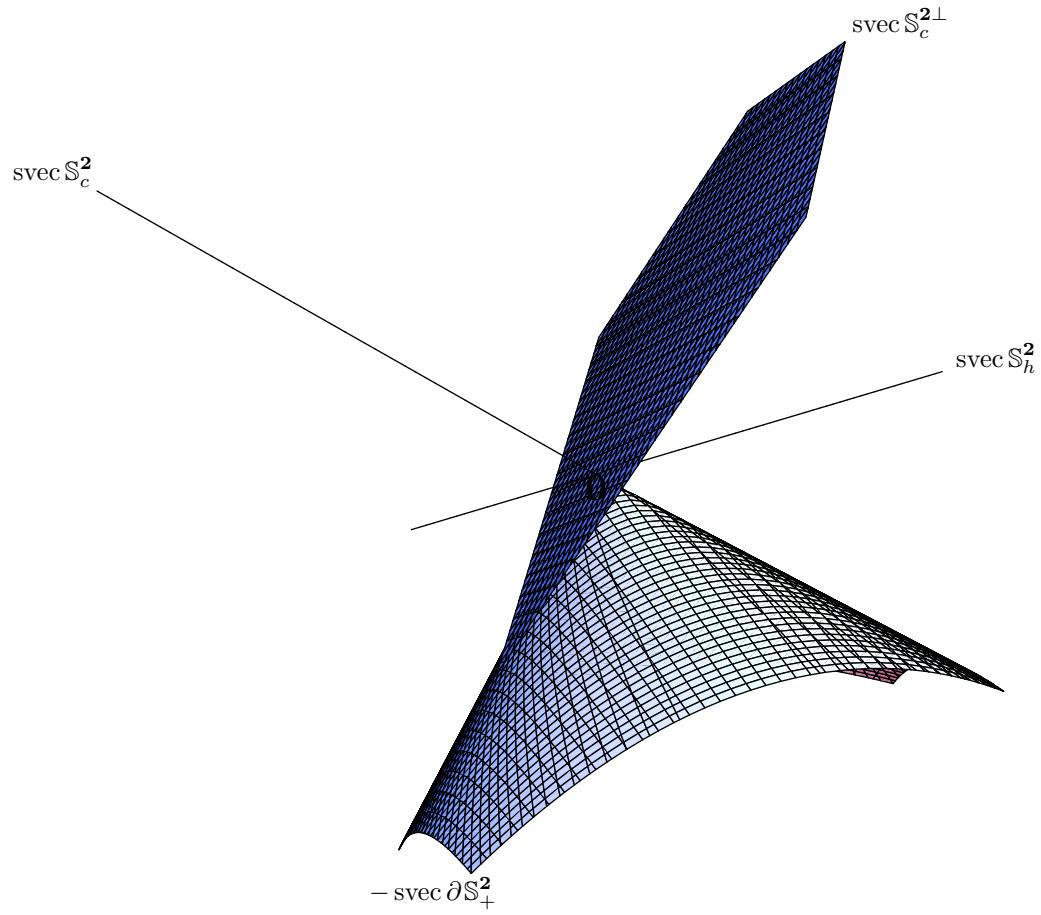


Figure 174: EDM cone construction in isometrically isomorphic  $\mathbb{R}^3$  by adding polar PSD cone to  $\text{svec } \mathbb{S}_c^{2\perp}$ . Difference  $\text{svec } (\mathbb{S}_c^{2\perp} - \mathbb{S}_+^2)$  is halfspace partially bounded by  $\text{svec } \mathbb{S}_c^{2\perp}$ . EDM cone is nonnegative halfline along  $\text{svec } \mathbb{S}_h^2$  in this dimension.

The dual EDM cone follows directly from (1414) by standard properties of cones (§2.13.1.2):

$$\text{EDM}^{N^*} = \overline{\mathcal{K}_1^* + \mathcal{K}_2^*} = \mathbb{S}_h^{N\perp} - \mathbb{S}_c^N \cap \mathbb{S}_+^N \quad (1415)$$

which bears strong resemblance to (1394).

### 6.8.1.2 nonnegative orthant contains $\text{EDM}^N$

That  $\text{EDM}^N$  is a proper subset of the nonnegative orthant is not obvious from (1414). We wish to verify

$$\text{EDM}^N = \mathbb{S}_h^N \cap (\mathbb{S}_c^{N\perp} - \mathbb{S}_+^N) \subset \mathbb{R}_+^{N \times N} \quad (1416)$$

While there are many ways to prove this, it is sufficient to show that all entries of the extreme directions of  $\text{EDM}^N$  must be nonnegative; *id est*, for any particular nonzero vector  $z = [z_i, i=1 \dots N] \in \mathcal{N}(\mathbf{1}^T)$  (§6.4.3.2),

$$\delta(zz^T)\mathbf{1}^T + \mathbf{1}\delta(zz^T)^T - 2zz^T \geq \mathbf{0} \quad (1417)$$

where the inequality denotes entrywise comparison. The inequality holds because the  $ij^{\text{th}}$  entry of an extreme direction is squared:  $(z_i - z_j)^2$ .

We observe that the dyad  $2zz^T \in \mathbb{S}_+^N$  belongs to the positive semidefinite cone, the doublet

$$\delta(zz^T)\mathbf{1}^T + \mathbf{1}\delta(zz^T)^T \in \mathbb{S}_c^{N\perp} \quad (1418)$$

to the orthogonal complement (2198) of the geometric center subspace, while their difference is a member of the symmetric hollow subspace  $\mathbb{S}_h^N$ . ♦

Here is an algebraic method to prove nonnegativity: Suppose we are given  $A \in \mathbb{S}_c^{N\perp}$  and  $B = [B_{ij}] \in \mathbb{S}_+^N$  and  $A - B \in \mathbb{S}_h^N$ . Then we have, for some vector  $u$ ,  $A = u\mathbf{1}^T + \mathbf{1}u^T = [A_{ij}] = [u_i + u_j]$  and  $\delta(B) = \delta(A) = 2u$ . Positive semidefiniteness of  $B$  requires nonnegativity  $A - B \geq \mathbf{0}$  because

$$(e_i - e_j)^T B (e_i - e_j) = (B_{ii} - B_{ij}) - (B_{ji} - B_{jj}) = 2(u_i + u_j) - 2B_{ij} \geq 0 \quad (1419)$$

♦

### 6.8.1.3 Dual Euclidean distance matrix criterion

Conditions necessary for membership of a matrix  $D^* \in \mathbb{S}^N$  to the dual EDM cone  $\text{EDM}^{N^*}$  may be derived from (1394):  $D^* \in \text{EDM}^{N^*} \Rightarrow D^* = \delta(y) - V_N A V_N^T$  for some vector  $y$  and positive semidefinite matrix  $A \in \mathbb{S}_+^{N-1}$ . This in turn implies  $\delta(D^* \mathbf{1}) = \delta(y)$ . Then, for  $D^* \in \mathbb{S}^N$

$$D^* \in \text{EDM}^{N^*} \Leftrightarrow \delta(D^* \mathbf{1}) - D^* \succeq 0 \quad (1420)$$

where, for any symmetric matrix  $D^*$

$$\delta(D^* \mathbf{1}) - D^* \in \mathbb{S}_c^N \quad (1421)$$

To show sufficiency of the matrix criterion in (1420), recall Gram-form EDM operator

$$\mathbf{D}(G) = \delta(G)\mathbf{1}^T + \mathbf{1}\delta(G)^T - 2G \quad (1045)$$

where Gram matrix  $G$  is positive semidefinite by definition, and recall the selfadjointness property of the main-diagonal linear operator  $\delta$  (§A.1):

$$\langle D, D^* \rangle = \langle \delta(G)\mathbf{1}^T + \mathbf{1}\delta(G)^T - 2G, D^* \rangle = \langle G, \delta(D^* \mathbf{1}) - D^* \rangle 2 \quad (1063)$$

Assuming  $\langle G, \delta(D^* \mathbf{1}) - D^* \rangle \geq 0$  (1639), then we have known membership relation (§2.13.2.0.1)

$$D^* \in \mathbb{EDM}^{N^*} \Leftrightarrow \langle D, D^* \rangle \geq 0 \quad \forall D \in \mathbb{EDM}^N \quad (1422)$$

◆

Elegance of this matrix criterion (1420) for membership to the dual EDM cone derives from lack of any other assumptions except that  $D^*$  be symmetric.<sup>6.14</sup> Linear Gram-form EDM operator  $\mathbf{D}(Y)$  (1045) has adjoint, for  $Y \in \mathbb{S}^N$

$$\mathbf{D}^T(Y) \triangleq (\delta(Y\mathbf{1}) - Y)2 \quad (1423)$$

Then from (1422) and (1046) we have: [97, p.111]

$$\begin{aligned} \mathbb{EDM}^{N^*} &= \{D^* \in \mathbb{S}^N \mid \langle D, D^* \rangle \geq 0 \quad \forall D \in \mathbb{EDM}^N\} \\ &= \{D^* \in \mathbb{S}^N \mid \langle \mathbf{D}(G), D^* \rangle \geq 0 \quad \forall G \in \mathbb{S}_+^N\} \\ &= \{D^* \in \mathbb{S}^N \mid \langle G, \mathbf{D}^T(D^*) \rangle \geq 0 \quad \forall G \in \mathbb{S}_+^N\} \\ &= \{D^* \in \mathbb{S}^N \mid \delta(D^* \mathbf{1}) - D^* \succeq 0\} \end{aligned} \quad (1424)$$

the dual EDM cone expressed in terms of the adjoint operator. A dual EDM cone determined this way is illustrated in Figure 176.

#### 6.8.1.3.1 Exercise. Dual EDM spectral cone.

Find a spectral cone as in §5.11.2 corresponding to  $\mathbb{EDM}^{N^*}$ . ▼

#### 6.8.1.4 Nonorthogonal components of dual EDM

Now we tie construct (1415) for the dual EDM cone together with the matrix criterion (1420) for dual EDM cone membership. For any  $D^* \in \mathbb{S}^N$  it is obvious:

$$\delta(D^* \mathbf{1}) \in \mathbb{S}_h^{N \perp} \quad (1425)$$

any diagonal matrix belongs to the subspace of diagonal matrices (68). We know when  $D^* \in \mathbb{EDM}^{N^*}$

$$\delta(D^* \mathbf{1}) - D^* \in \mathbb{S}_c^N \cap \mathbb{S}_+^N \quad (1426)$$

this adjoint expression (1423) belongs to that face (1384) of the positive semidefinite cone  $\mathbb{S}_+^N$  in the geometric center subspace. Any nonzero dual EDM

$$D^* = \delta(D^* \mathbf{1}) - (\delta(D^* \mathbf{1}) - D^*) \in \mathbb{S}_h^{N \perp} \ominus \mathbb{S}_c^N \cap \mathbb{S}_+^N = \mathbb{EDM}^{N^*} \quad (1427)$$

can therefore be expressed as the difference of two linearly independent (when vectorized) nonorthogonal components (Figure 154, Figure 175).

#### 6.8.1.5 Affine dimension complementarity

From §6.8.1.3 we have, for some  $A \in \mathbb{S}_+^{N-1}$  (confer (1426))

$$\delta(D^* \mathbf{1}) - D^* = V_{\mathcal{N}} A V_{\mathcal{N}}^T \in \mathbb{S}_c^N \cap \mathbb{S}_+^N \quad (1428)$$

if and only if  $D^*$  belongs to the dual EDM cone. Call  $\text{rank}(V_{\mathcal{N}} A V_{\mathcal{N}}^T)$  *dual affine dimension*. Empirically, we find a complementary relationship in affine dimension

---

<sup>6.14</sup>Recall: Schoenberg criterion (1052), for membership to the EDM cone, requires membership to the symmetric hollow subspace.

$$D^\circ = \delta(D^\circ \mathbf{1}) + (D^\circ - \delta(D^\circ \mathbf{1})) \in \mathbb{S}_h^{N\perp} \oplus \mathbb{S}_c^N \cap \mathbb{S}_+^N = \text{EDM}^{N^\circ}$$

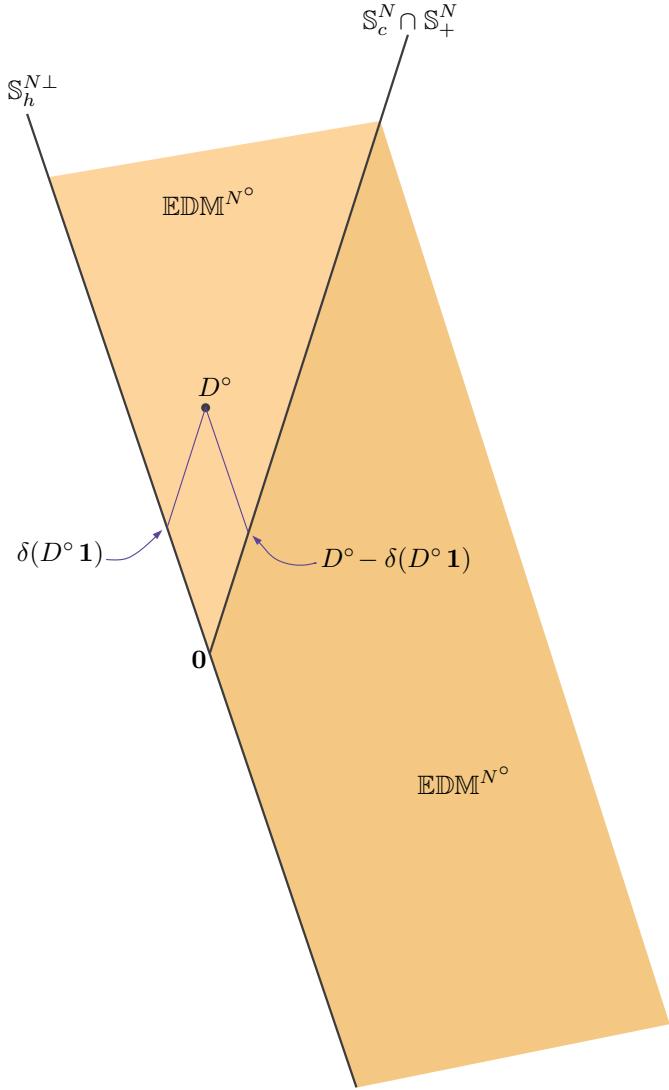


Figure 175: [Hand-drawn abstraction](#) of polar EDM cone  $\text{EDM}^{N^\circ}$  (drawn truncated). Any member  $D^\circ$  of polar EDM cone can be decomposed into two linearly independent nonorthogonal components:  $\delta(D^\circ \mathbf{1})$  and  $D^\circ - \delta(D^\circ \mathbf{1})$ .

between the projection of some arbitrary symmetric matrix  $H$  on the polar EDM cone,  $\mathbb{EDM}^{N^\circ} = -\mathbb{EDM}^{N^*}$ , and its projection on the EDM cone; *id est*, the optimal solution of 6.15

$$\begin{aligned} & \underset{D^\circ \in \mathbb{S}_h^N}{\text{minimize}} \quad \|D^\circ - H\|_F \\ & \text{subject to} \quad D^\circ - \delta(D^\circ \mathbf{1}) \succeq 0 \end{aligned} \quad (1429)$$

has dual affine dimension complementary to affine dimension corresponding to the optimal solution of

$$\begin{aligned} & \underset{D \in \mathbb{S}_h^N}{\text{minimize}} \quad \|D - H\|_F \\ & \text{subject to} \quad -V_N^T D V_N \succeq 0 \end{aligned} \quad (1430)$$

Precisely,

$$\text{rank}(D^{\circ*} - \delta(D^{\circ*} \mathbf{1})) + \text{rank}(V_N^T D^* V_N) = N - 1 \quad (1431)$$

and  $\text{rank}(D^{\circ*} - \delta(D^{\circ*} \mathbf{1})) \leq N - 1$  because vector  $\mathbf{1}$  is always in the nullspace of rank's argument. This is similar to the known result for projection on the selfdual positive semidefinite cone and its polar:

$$\text{rank } P_{\mathbb{S}_+^N} H + \text{rank } P_{\mathbb{S}_+^N} H = N \quad (1432)$$

When low affine dimension is a desirable result of projection on the EDM cone, projection on the polar EDM cone should be performed instead. Convex polar problem (1429) can be solved for  $D^{\circ*}$  by transforming to an equivalent Schur-form semidefinite program (§3.5.3). Interior-point methods, for numerically solving semidefinite programs, tend to produce high-rank solutions. (§4.1.2) Then  $D^* = H - D^{\circ*} \in \mathbb{EDM}^N$  by Corollary E.9.2.2.1, and  $D^*$  will tend to have low affine dimension. This approach breaks when attempting projection on a cone subset discriminated by affine dimension or rank, because then we have no complementarity relation like (1431) or (1432) (§7.1.4.1).

### 6.8.1.6 EDM cone is not selfdual

In §5.6.1.1, via Gram-form EDM operator

$$\mathbf{D}(G) = \delta(G)\mathbf{1}^T + \mathbf{1}\delta(G)^T - 2G \in \mathbb{EDM}^N \iff G \succeq 0 \quad (1045)$$

we established clear connection between the EDM cone and that face (1384) of positive semidefinite cone  $\mathbb{S}_+^N$  in the geometric center subspace:

$$\mathbb{EDM}^N = \mathbf{D}(\mathbb{S}_c^N \cap \mathbb{S}_+^N) \quad (1152)$$

$$\mathbf{V}(\mathbb{EDM}^N) = \mathbb{S}_c^N \cap \mathbb{S}_+^N \quad (1153)$$

where

$$\mathbf{V}(D) = -VDV^{\frac{1}{2}} \quad (1141)$$

In §5.6.1 we established

$$\mathbb{S}_c^N \cap \mathbb{S}_+^N = V_N \mathbb{S}_+^{N-1} V_N^T \quad (1139)$$

---

<sup>6.15</sup>This polar projection can be solved quickly (without semidefinite programming) via Lemma 6.8.1.1.1; rewriting,

$$\begin{aligned} & \underset{D^\circ \in \mathbb{S}_h^N}{\text{minimize}} \quad \|(D^\circ - \delta(D^\circ \mathbf{1})) - (H - \delta(D^\circ \mathbf{1}))\|_F \\ & \text{subject to} \quad D^\circ - \delta(D^\circ \mathbf{1}) \succeq 0 \end{aligned}$$

which is the projection of affinely transformed optimal solution  $H - \delta(D^{\circ*} \mathbf{1})$  on  $\mathbb{S}_c^N \cap \mathbb{S}_+^N$ ;

$$D^{\circ*} - \delta(D^{\circ*} \mathbf{1}) = P_{\mathbb{S}_+^N} P_{\mathbb{S}_c^N} (H - \delta(D^{\circ*} \mathbf{1}))$$

Foreknowledge of an optimal solution  $D^{\circ*}$  as argument to projection suggests recursion.

Then from (1420), (1428), and (1394) we can deduce

$$\delta(\mathbb{EDM}^{N^*} \mathbf{1}) - \mathbb{EDM}^{N^*} = V_{\mathcal{N}} \mathbb{S}_+^{N-1} V_{\mathcal{N}}^T = \mathbb{S}_c^N \cap \mathbb{S}_+^N \quad (1433)$$

which, by (1152) and (1153), means the EDM cone can be related to the dual EDM cone by an equality:

$$\mathbb{EDM}^N = \mathbf{D}(\delta(\mathbb{EDM}^{N^*} \mathbf{1}) - \mathbb{EDM}^{N^*}) \quad (1434)$$

$$\mathbf{V}(\mathbb{EDM}^N) = \delta(\mathbb{EDM}^{N^*} \mathbf{1}) - \mathbb{EDM}^{N^*} \quad (1435)$$

This means projection  $-\mathbf{V}(\mathbb{EDM}^N)$  of the EDM cone on the geometric center subspace  $\mathbb{S}_c^N$  (§E.7.2.0.2) is a linear transformation of the dual EDM cone:  $\mathbb{EDM}^{N^*} - \delta(\mathbb{EDM}^{N^*} \mathbf{1})$ . Secondarily, it means the EDM cone is not selfdual in  $\mathbb{S}^N$ .

#### 6.8.1.7 Schoenberg criterion is discretized membership relation

We show the Schoenberg criterion

$$\left. \begin{array}{l} -V_{\mathcal{N}}^T D V_{\mathcal{N}} \in \mathbb{S}_+^{N-1} \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1052)$$

to be a discretized membership relation (§2.13.4) between a closed convex cone  $\mathcal{K}$  and its dual  $\mathcal{K}^*$  like

$$\langle y, x \rangle \geq 0 \text{ for all } y \in \mathcal{G}(\mathcal{K}^*) \Leftrightarrow x \in \mathcal{K} \quad (369)$$

where  $\mathcal{G}(\mathcal{K}^*)$  is any set of generators whose conic hull constructs closed convex dual cone  $\mathcal{K}^*$ :

The Schoenberg criterion is the same as

$$\left. \begin{array}{l} \langle zz^T, -D \rangle \geq 0 \quad \forall zz^T \mid \mathbf{1}\mathbf{1}^T zz^T = \mathbf{0} \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1378)$$

which, by (1379), is the same as

$$\left. \begin{array}{l} \langle zz^T, -D \rangle \geq 0 \quad \forall zz^T \in \left\{ V_{\mathcal{N}} v v^T V_{\mathcal{N}}^T \mid v \in \mathbb{R}^{N-1} \right\} \\ D \in \mathbb{S}_h^N \end{array} \right\} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1436)$$

where the  $zz^T$  constitute a set of generators  $\mathcal{G}$  for the positive semidefinite cone's smallest face  $\mathcal{F}(\mathbb{S}_+^N \ni V)$  (§6.6.1) that contains auxiliary matrix  $V$ . From the aggregate in (1394) we get the ordinary membership relation, assuming only  $D \in \mathbb{S}^N$  [225, p.58]

$$\langle D^*, D \rangle \geq 0 \quad \forall D^* \in \mathbb{EDM}^{N^*} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1437)$$

$$\langle D^*, D \rangle \geq 0 \quad \forall D^* \in \{\delta(u) \mid u \in \mathbb{R}^N\} - \text{cone}\left\{ V_{\mathcal{N}} v v^T V_{\mathcal{N}}^T \mid v \in \mathbb{R}^{N-1} \right\} \Leftrightarrow D \in \mathbb{EDM}^N$$

Discretization (369) yields:

$$\langle D^*, D \rangle \geq 0 \quad \forall D^* \in \{e_i e_i^T, -e_j e_j^T, -V_{\mathcal{N}} v v^T V_{\mathcal{N}}^T \mid i, j = 1 \dots N, v \in \mathbb{R}^{N-1}\} \Leftrightarrow D \in \mathbb{EDM}^N \quad (1438)$$

Because  $\langle \{\delta(u) \mid u \in \mathbb{R}^N\}, D \rangle \geq 0 \Leftrightarrow D \in \mathbb{S}_h^N$ , we can restrict observation to the symmetric hollow subspace without loss of generality. Then for  $D \in \mathbb{S}_h^N$

$$\langle D^*, D \rangle \geq 0 \quad \forall D^* \in \left\{ -V_N v v^T V_N^T \mid v \in \mathbb{R}^{N-1} \right\} \Leftrightarrow D \in \text{EDM}^N \quad (1439)$$

this discretized membership relation becomes (1436); identical to the Schoenberg criterion.

Hitherto a correspondence between the EDM cone and a face of a PSD cone, the Schoenberg criterion is now accurately interpreted as a discretized membership relation between the EDM cone and its ordinary dual.

### 6.8.2 Ambient $\mathbb{S}_h^N$

When instead we consider the ambient space of symmetric hollow matrices (1395), then still we find the EDM cone is not selfdual for  $N > 2$ . The simplest way to prove this is as follows:

Given a set of generators  $\mathcal{G} = \{\Gamma\}$  (1355) for the pointed closed convex EDM cone, the *discretized membership theorem* in §2.13.4.2.1 asserts that members of the dual EDM cone in the ambient space of symmetric hollow matrices can be discerned via discretized membership relation:

$$\begin{aligned} \text{EDM}^{N^*} \cap \mathbb{S}_h^N &\triangleq \{D^* \in \mathbb{S}_h^N \mid \langle \Gamma, D^* \rangle \geq 0 \quad \forall \Gamma \in \mathcal{G}(\text{EDM}^N)\} \\ &= \{D^* \in \mathbb{S}_h^N \mid \langle \delta(z z^T) \mathbf{1}^T + \mathbf{1} \delta(z z^T)^T - 2 z z^T, D^* \rangle \geq 0 \quad \forall z \in \mathcal{N}(\mathbf{1}^T)\} \\ &= \{D^* \in \mathbb{S}_h^N \mid \langle \mathbf{1} \delta(z z^T)^T - z z^T, D^* \rangle \geq 0 \quad \forall z \in \mathcal{N}(\mathbf{1}^T)\} \end{aligned} \quad (1440)$$

By comparison

$$\text{EDM}^N = \{D \in \mathbb{S}_h^N \mid \langle -z z^T, D \rangle \geq 0 \quad \forall z \in \mathcal{N}(\mathbf{1}^T)\} \quad (1441)$$

the term  $\delta(z z^T)^T D^* \mathbf{1}$  foils any hope of selfduality in ambient  $\mathbb{S}_h^N$ . ♦

To find the dual EDM cone in ambient  $\mathbb{S}_h^N$  per §2.13.10.4 we prune the aggregate in (1394) describing the ordinary dual EDM cone, removing any member having nonzero main diagonal:

$$\begin{aligned} \text{EDM}^{N^*} \cap \mathbb{S}_h^N &= \text{cone} \left\{ \delta^2(V_N v v^T V_N^T) - V_N v v^T V_N^T \mid v \in \mathbb{R}^{N-1} \right\} \\ &= \{ \delta^2(V_N \Psi V_N^T) - V_N \Psi V_N^T \mid \Psi \in \mathbb{S}_+^{N-1} \} \end{aligned} \quad (1442)$$

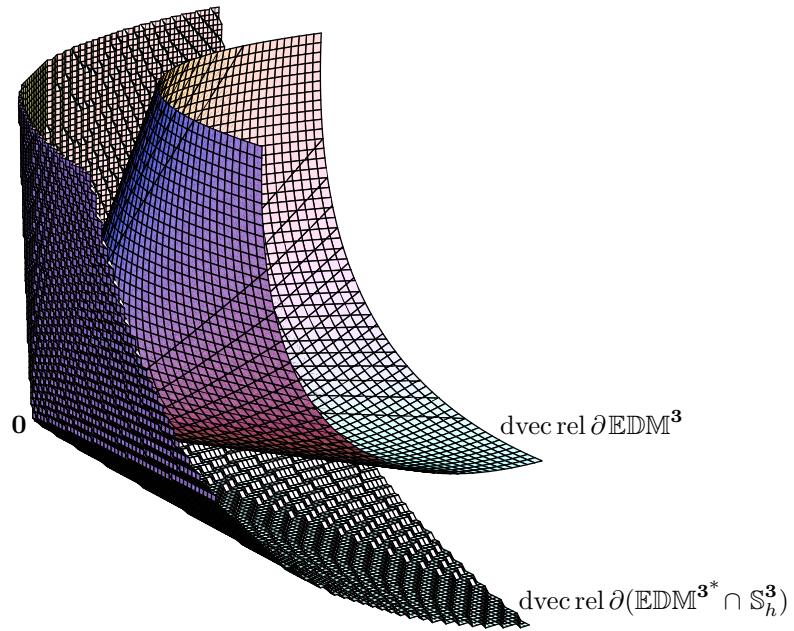
When  $N = 1$ , the EDM cone and its dual in ambient  $\mathbb{S}_h$  each comprise the origin in isomorphic  $\mathbb{R}^0$ ; thus, selfdual in this dimension. (confer(106))

When  $N = 2$ , the EDM cone is the nonnegative real line in isomorphic  $\mathbb{R}$ . (Figure 167)  $\text{EDM}^{2^*}$  in  $\mathbb{S}_h^2$  is identical, thus selfdual in this dimension. This result is in agreement with (1440), verified directly: for all  $\kappa \in \mathbb{R}$ ,  $z = \kappa \begin{bmatrix} 1 \\ -1 \end{bmatrix}$  and  $\delta(z z^T) = \kappa^2 \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow d_{12}^* \geq 0$ .

The first case adverse to selfduality  $N = 3$  may be deduced from Figure 163; the EDM cone is a circular cone in isomorphic  $\mathbb{R}^3$  corresponding to no rotation of Lorentz cone (181) (the selfdual circular cone). Figure 176 illustrates the EDM cone and its dual in ambient  $\mathbb{S}_h^3$ ; no longer selfdual.

#### 6.8.2.0.1 Exercise. Positive semidefinite cone from EDM cone.

What, if any, is the inversion of semidefinite and distance cone equality (1414)? That is to say, can  $\mathbb{S}_+$  be expressed only in terms of  $\text{EDM}$ ,  $\mathbb{S}_h$ , and  $\mathbb{S}_c$ ? ▼



$$D^* \in \text{EDM}^{N^*} \Leftrightarrow \delta(D^* \mathbf{1}) - D^* \succeq 0 \quad (\text{1420})$$

Figure 176: Ordinary dual EDM cone projected on  $S_h^3$  shrouds  $\text{EDM}^3$ ; drawn tiled in isometrically isomorphic  $\mathbb{R}^3$ . (It so happens: intersection  $\text{EDM}^{N^*} \cap S_h^N$  (§2.13.10.3) is identical to projection of dual EDM cone on  $S_h^N$ .)

**6.8.2.0.2 Exercise.** *Rank complementarity for EDM cone.*  
Prove (1431). ▼

## 6.9 Theorem of the alternative

In §2.13.2.1.1 we showed how alternative systems of generalized inequality can be derived from closed convex cones and their duals. This section is, therefore, a fitting postscript to the discussion of the dual EDM cone.

**6.9.0.0.1 Theorem.** *EDM alternative.*

[187, §1]

Given  $D \in \mathbb{S}_h^N$

$$D \in \text{EDM}^N$$

or in the alternative

$$\exists z \text{ such that } \begin{cases} \mathbf{1}^T z = 1 \\ Dz = \mathbf{0} \end{cases} \quad (1443)$$

In words, either  $\mathcal{N}(D)$  intersects hyperplane  $\{z \mid \mathbf{1}^T z = 1\}$  or  $D$  is an EDM; the alternatives are incompatible. ◇

When  $D$  is an EDM [290, §2]

$$\mathcal{N}(D) \subset \mathcal{N}(\mathbf{1}^T) = \{z \mid \mathbf{1}^T z = 0\} \quad (1444)$$

Because [187, §2] (§E.0.1)

$$\begin{aligned} DD^\dagger \mathbf{1} &= \mathbf{1} \\ \mathbf{1}^T D^\dagger D &= \mathbf{1}^T \end{aligned} \quad (1445)$$

then

$$\mathcal{R}(\mathbf{1}) \subset \mathcal{R}(D) \quad (1446)$$

## 6.10 Postscript

We provided an equality (1414) relating the convex cone of Euclidean distance matrices to the convex cone of positive semidefinite matrices. Projection on a positive semidefinite cone, constrained by an upper bound on rank, is easy and well known; [147] simply, a matter of truncating a list of eigenvalues. Projection on a positive semidefinite cone with such a rank constraint is, in fact, a convex optimization problem. (§7.1.4)

In the past, it was difficult to project on the EDM cone under a constraint on rank or affine dimension. A surrogate method was to invoke the Schoenberg criterion (1052) and then project on a positive semidefinite cone under a rank constraint bounding affine dimension from above. But a solution acquired that way is necessarily suboptimal.

In §7.3.3 we present a method for projecting directly on the EDM cone under a constraint on rank or affine dimension.



# Chapter 7

## Proximity problems

*In the “extremely large-scale case” ( $N$  of order of tens and hundreds of thousands), [iteration cost  $O(N^3)$ ] rules out all advanced convex optimization techniques, including all known polynomial time algorithms.*

— Arkadi Nemirovski, 2004

A problem common to various sciences is to find the Euclidean distance matrix (EDM)  $D \in \mathbb{EDM}^N$  closest in some sense to a given complete matrix of measurements  $H$  under a constraint on affine dimension  $0 \leq r \leq N - 1$  (§2.3.1, §5.7.1.1); rather,  $r$  is bounded above by desired affine dimension  $\rho$ .

### 7.0.1 Measurement matrix $H$

Ideally, we want a given matrix of measurements  $H \in \mathbb{R}^{N \times N}$  to conform with the first three Euclidean metric properties (§5.2); to belong to the intersection of the orthant of nonnegative matrices  $\mathbb{R}_+^{N \times N}$  with the symmetric hollow subspace  $\mathbb{S}_h^N$  (§2.2.3.0.1). Geometrically, we want  $H$  to belong to the polyhedral cone (§2.12.1.0.1)

$$\mathcal{K} \triangleq \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N} \quad (1447)$$

Yet in practice,  $H$  can possess significant measurement uncertainty (noise).

Sometimes realization of an optimization problem demands that its input, the given matrix  $H$ , possess some particular characteristics; perhaps symmetry and hollowness or nonnegativity. When that  $H$  given does not have the desired properties, then we must impose them upon  $H$  prior to optimization:

- When *measurement matrix*  $H$  is neither symmetric or hollow, taking its symmetric hollow part is equivalent to orthogonal projection on the symmetric hollow subspace  $\mathbb{S}_h^N$ .
- When measurements of distance in  $H$  are negative, zeroing negative entries effects unique minimum-distance projection on the orthant of nonnegative matrices  $\mathbb{R}_+^{N \times N}$  in isomorphic  $\mathbb{R}^{N^2}$  (§E.9.2.2.3).

### 7.0.1.1 Order of imposition

Since convex cone  $\mathcal{K}$  (1447) is the intersection of an orthant with a subspace, we want to project on that subset of the orthant belonging to the subspace; on the nonnegative orthant in the symmetric hollow subspace that is, in fact, the intersection. For that reason alone, unique minimum-distance projection of  $H$  on  $\mathcal{K}$  (that member of  $\mathcal{K}$  closest to  $H$  in isomorphic  $\mathbb{R}^{N^2}$  in the Euclidean sense) can be attained by first taking its symmetric hollow part, and only then clipping negative entries of the result to 0; *id est*, there is only one correct *order of projection*, in general, on an orthant intersecting a subspace:

- project on the subspace, then project the result on the orthant in that subspace. (*confer* §E.9.5)

In contrast, order of projection on an intersection of subspaces is arbitrary.

That order of projection rule applies more generally, of course, to intersection of any convex set  $\mathcal{C}$  with any subspace. Consider the *proximity problem*<sup>7.1</sup> over convex feasible set  $\mathbb{S}_h^N \cap \mathcal{C}$  given nonsymmetric nonhollow  $H \in \mathbb{R}^{N \times N}$ :

$$\begin{aligned} & \underset{B \in \mathbb{S}_h^N}{\text{minimize}} \quad \|B - H\|_F^2 \\ & \text{subject to} \quad B \in \mathcal{C} \end{aligned} \tag{1448}$$

a convex optimization problem. Because the symmetric hollow subspace  $\mathbb{S}_h^N$  is orthogonal to the antisymmetric antihollow subspace  $\mathbb{R}_h^{N \times N \perp}$  (§2.2.3), then for  $B \in \mathbb{S}_h^N$

$$\text{tr}\left(B^T \left(\frac{1}{2}(H - H^T) + \delta^2(H)\right)\right) = 0 \tag{1449}$$

so the objective function is equivalent to

$$\|B - H\|_F^2 \equiv \left\|B - \left(\frac{1}{2}(H + H^T) - \delta^2(H)\right)\right\|_F^2 + \left\|\frac{1}{2}(H - H^T) + \delta^2(H)\right\|_F^2 \tag{1450}$$

This means the antisymmetric antihollow part of given matrix  $H$  would be ignored by minimization with respect to symmetric hollow variable  $B$  under Frobenius' norm; *id est*, minimization proceeds as though given the symmetric hollow part of  $H$ .

This action of Frobenius' norm (1450) is effectively a Euclidean projection (minimum-distance projection) of  $H$  on the symmetric hollow subspace  $\mathbb{S}_h^N$  prior to minimization. Thus minimization inherently follows the correct order for projection on  $\mathbb{S}_h^N \cap \mathcal{C}$ . Therefore we may either assume  $H \in \mathbb{S}_h^N$ , or take its symmetric hollow part prior to optimization.

### 7.0.1.2 Flagrant input error under nonnegativity demand

More pertinent to the optimization problems presented herein where

$$\mathcal{C} \triangleq \mathbb{EDM}^N \subseteq \mathcal{K} = \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N} \tag{1451}$$

then should some particular realization of a proximity problem demand input  $H$  be nonnegative, and were we only to zero negative entries of a nonsymmetric nonhollow input  $H$  prior to optimization, then the ensuing projection on  $\mathbb{EDM}^N$  would be guaranteed incorrect (out of order).

---

<sup>7.1</sup>There are two equivalent interpretations of projection (§E.9): one finds a set normal, the other, minimum distance between a point and a set. Here we realize the latter view.

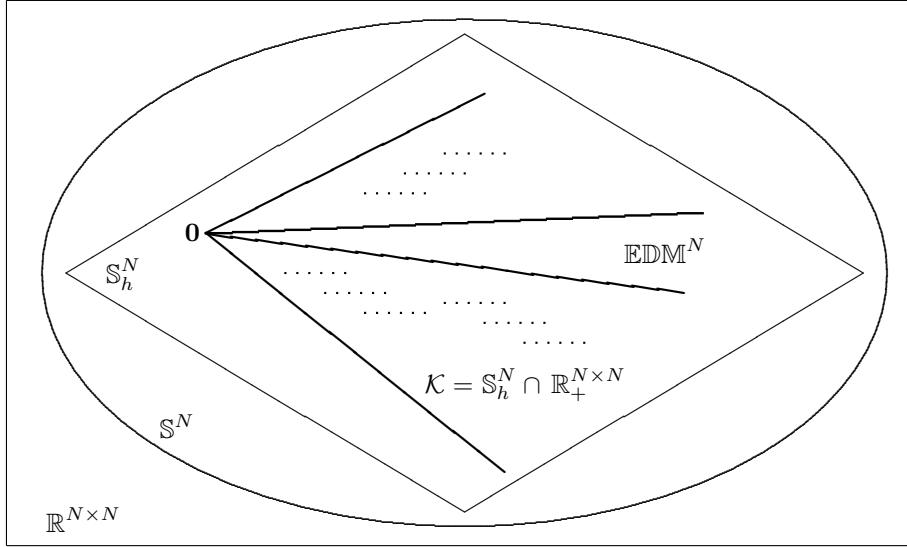


Figure 177: Pseudo-Venn diagram: EDM cone  $\text{EDM}^N$  belongs to intersection of symmetric hollow subspace with nonnegative orthant;  $\text{EDM}^N \subseteq \mathcal{K}$  (1032).  $\text{EDM}^N$  cannot exist outside of  $\mathbb{S}_h^N$ , but  $\mathbb{R}_+^{N \times N}$  does.

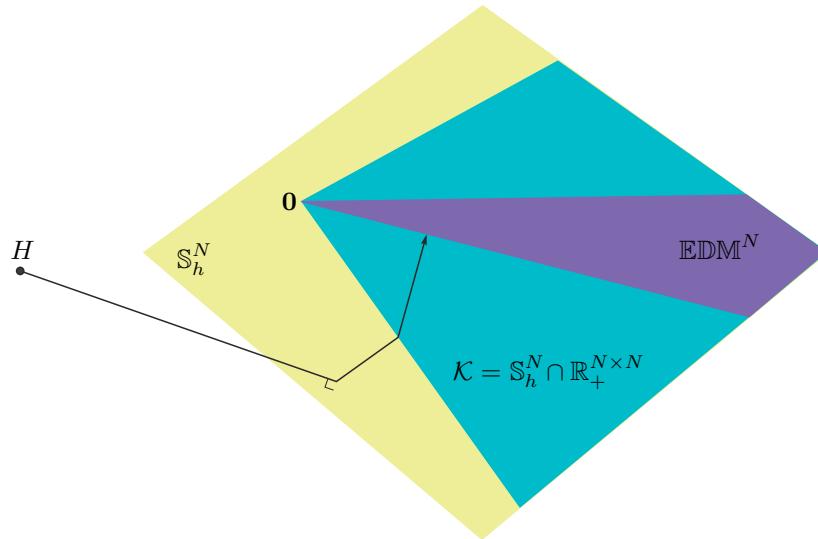


Figure 178: Pseudo-Venn diagram from Figure 177 showing elbow placed in path of projection of  $H$  on  $\text{EDM}^N \subset \mathbb{S}_h^N$  by an optimization problem demanding nonnegative input matrix  $H$ . The first two line segments, leading away from  $H$ , result from correct order of projection required to provide nonnegative  $H$  prior to optimization. Were  $H$  nonnegative, then its projection on  $\mathbb{S}_h^N$  would instead belong to  $\mathcal{K}$ ; making the elbow disappear. (confer Figure 194)

Now comes a surprising fact: Even were we to correctly follow the order of projection rule so as to provide  $H \in \mathcal{K}$  prior to optimization, then the ensuing projection on  $\text{EDM}^N$  will be incorrect whenever input  $H$  has negative entries and some proximity problem demands nonnegative input  $H$ .

This is best understood referring to Figure 177: Suppose nonnegative input  $H$  is demanded, and then the problem realization correctly projects its input first on  $\mathbb{S}_h^N$  and then directly on  $\mathcal{C} = \text{EDM}^N$ . That demand for nonnegativity effectively requires imposition of  $\mathcal{K}$  on input  $H$  prior to optimization so as to obtain correct order of projection (on  $\mathbb{S}_h^N$  first). Yet such an imposition prior to projection on  $\text{EDM}^N$  generally introduces an *elbow* into the path of projection (illustrated in Figure 178) caused by the technique itself; that being, a particular proximity problem realization requiring nonnegative input.

Any procedure, for imposition of nonnegativity on input  $H$ , can only be incorrect in this circumstance. There is no resolution unless input  $H$  is guaranteed nonnegative with no tinkering. Otherwise, we have no choice but to employ a different problem realization; one not demanding nonnegative input.

### 7.0.2 Least lower bound

Most of the problems we encounter in this chapter have the general form:

$$\begin{aligned} & \underset{B}{\text{minimize}} && \|B - A\|_{\text{F}} \\ & \text{subject to} && B \in \mathcal{C} \end{aligned} \tag{1452}$$

where  $A \in \mathbb{R}^{m \times n}$  is given data. This particular objective denotes Euclidean projection (§E) of vectorized matrix  $A$  on the set  $\mathcal{C}$  which may or may not be convex. When  $\mathcal{C}$  is convex, then projection is unique minimum-distance because Frobenius' norm square is a strictly convex function of variable  $B$  and because the optimal solution is the same regardless of the square (518). When  $\mathcal{C}$  is a subspace, then the direction of projection is orthogonal to  $\mathcal{C}$ .

Denoting by  $A \triangleq U_A \Sigma_A Q_A^T$  and  $B \triangleq U_B \Sigma_B Q_B^T$  their full singular value decompositions (whose singular values are always nonincreasingly ordered (§A.6)), there exists a tight lower bound on the objective over the manifold of orthogonal matrices;

$$\|\Sigma_B - \Sigma_A\|_{\text{F}} \leq \inf_{U_A, U_B, Q_A, Q_B} \|B - A\|_{\text{F}} \tag{1453}$$

This least lower bound holds more generally for any orthogonally invariant norm on  $\mathbb{R}^{m \times n}$  (§2.2.1) including the Frobenius and spectral norm [364, §II.3]. [228, §7.4.51]

### 7.0.3 Problem approach.

*stress/sstress* problems traditionally posed in terms of point position  $\{x_i \in \mathbb{R}^n, i=1 \dots N\}$

$$\underset{\{x_i\}}{\text{minimize}} \sum_{i, j \in \mathcal{I}} (\|x_i - x_j\| - h_{ij})^2 \tag{1454}$$

$$\underset{\{x_i\}}{\text{minimize}} \sum_{i, j \in \mathcal{I}} (\|x_i - x_j\|^2 - h_{ij})^2 \tag{1455}$$

(where  $\mathcal{I}$  is an abstract set of indices and  $h_{ij}$  is given data) are everywhere converted herein to the distance-square variable  $D$  or to Gram matrix  $G$ ; the Gram matrix acting as bridge between position and distance. (That conversion is performed regardless of whether known data is complete.) Then the techniques of chapter 5 or chapter 6 are applied to find relative or absolute position. This approach is taken because we prefer introduction of rank constraints into convex problems rather than searching a googol of local minima in nonconvex problems like (1455) or (1454) [117] (§3.9.0.0.3, §7.2.2.7.1).

### 7.0.4 Three prevalent proximity problems

There are three statements of the closest-EDM problem prevalent in the literature, the multiplicity due primarily to choice of projection on the EDM *versus* positive semidefinite (PSD) cone and vacillation between the distance-square variable  $d_{ij}$  *versus* absolute distance  $\sqrt{d_{ij}}$ . In their most fundamental form, the three prevalent proximity problems are (1456.1), (1456.2), and (1456.3): [380] for  $D \triangleq [d_{ij}]$  and  $\sqrt[3]{D} \triangleq [\sqrt{d_{ij}}]$

$$\begin{aligned}
(1) \quad & \underset{D}{\text{minimize}} \quad \| -V(D - H)V \|^2_F & & \underset{\sqrt[3]{D}}{\text{minimize}} \quad \| \sqrt[3]{D} - H \|^2_F \\
& \text{subject to} \quad \text{rank } VDV \leq \rho & & \text{subject to} \quad \text{rank } VDV \leq \rho \\
& D \in \mathbb{EDM}^N & & \sqrt[3]{D} \in \sqrt{\mathbb{EDM}^N} \\
(2) \quad & & & (1456) \\
(3) \quad & \underset{D}{\text{minimize}} \quad \| D - H \|^2_F & & \underset{\sqrt[3]{D}}{\text{minimize}} \quad \| -V(\sqrt[3]{D} - H)V \|^2_F \\
& \text{subject to} \quad \text{rank } VDV \leq \rho & & \text{subject to} \quad \text{rank } VDV \leq \rho \\
& D \in \mathbb{EDM}^N & & \sqrt[3]{D} \in \sqrt{\mathbb{EDM}^N} \\
(4) \quad & & & (1456)
\end{aligned}$$

where we have made explicit an imposed upper bound  $\rho$  on affine dimension

$$r = \text{rank } V_N^T DV_N = \text{rank } VDV \quad (1187)$$

that is benign when  $\rho = N - 1$  or  $H$  were realizable with  $r \leq \rho$ . Problems (1456.2) and (1456.3) are Euclidean projections of vectorized matrix  $H$  on an EDM cone, whereas problems (1456.1) and (1456.4) are Euclidean projections of vectorized matrix  $-VHV$  on a PSD cone.<sup>7.2</sup> (§6.3) Problem (1456.4) is not posed in the literature because it has limited theoretical foundation.<sup>7.3</sup>

Analytical solution to (1456.1) is known in closed form for any bound  $\rho$  and any auxiliary matrix  $V$  although, as the problem is stated, it is a convex optimization only in the case  $\rho = N - 1$ . We show, in §7.1.4, how (1456.1) becomes a convex optimization problem for any  $\rho$  when transformed to the spectral domain. When expressed as a function of point list in a matrix  $X$  as in (1454), problem (1456.2) becomes a variant of what is known in statistics literature as a *stress problem*. [56, p.34] [115] [395] Problem (1456.3) is a rank-constrained *sstress problem*, whereas (1456.1) is equivalent to a rank-constrained *strain problem*. [116, §5]<sup>7.4</sup> Problems (1456.2) and (1456.3) are convex optimization problems in  $D$  for the case  $\rho = N - 1$  wherein (1456.3) becomes equivalent to (1455). Even with the rank constraint removed from (1456.2), we will see that the convex problem remaining inherently minimizes affine dimension.

Generally speaking, each problem in (1456) produces a different result because there is no isometry relating them. Of the various auxiliary  $V$ -matrices (§B.4), the geometric centering matrix  $V$  (1055) appears in the literature most often although  $V_N$  (1039) is the auxiliary matrix naturally consequent to Schoenberg's seminal exposition [349]. Substitution of any auxiliary matrix or its pseudoinverse into these problems produces another valid problem.

Substitution of  $V_N$  for  $V$  in (1456.1), in particular, produces a different result because

$$\begin{aligned}
& \underset{D}{\text{minimize}} \quad \| -V_N^T(D - H)V_N \|^2_F \\
& \text{subject to} \quad D \in \mathbb{EDM}^N
\end{aligned} \quad (1457)$$

<sup>7.2</sup>Because  $-VHV$  is orthogonal projection of  $-H$  on the geometric center subspace  $\mathbb{S}_c^N$  (§E.7.2.0.2), problems (1456.1) and (1456.4) may be interpreted as oblique (nonminimum distance) projections of  $-H$  on a positive semidefinite cone.

<sup>7.3</sup> $D \in \mathbb{EDM}^N \Rightarrow \sqrt[3]{D} \in \mathbb{EDM}^N, -V\sqrt[3]{D}V \in \mathbb{S}_+^N$  (§5.10)

<sup>7.4</sup>Equivalence to de Leeuw's strain problem statement is established for  $\rho = N - 1$  via (1806) (39) (44).

finds  $D$  to attain Euclidean distance of vectorized  $-V_N^T H V_N$  to the positive semidefinite cone in isometrically isomorphic subspace  $\mathbb{R}^{N(N-1)/2}$ , whereas

$$\begin{aligned} & \underset{D}{\text{minimize}} \quad \| -V(D - H)V \|^2_F \\ & \text{subject to} \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1458)$$

attains Euclidean distance of vectorized  $-V H V$  to the positive semidefinite cone in ambient isometrically isomorphic  $\mathbb{R}^{N(N+1)/2}$ ; quite different projections<sup>7.5</sup> regardless of whether affine dimension is constrained. But substitution of auxiliary matrix  $V_W^T$  (§B.4.3) or  $V_N^\dagger$  yields the same result as (1456.1) because  $V = V_W V_W^T = V_N V_N^\dagger$ ; *id est*,

$$\begin{aligned} \| -V(D - H)V \|^2_F &= \| -V_W V_W^T (D - H) V_W V_W^T \|^2_F = \| -V_W^T (D - H) V_W \|^2_F \\ &= \| -V_N V_N^\dagger (D - H) V_N V_N^\dagger \|^2_F = \| -V_N^\dagger (D - H) V_N \|^2_F \end{aligned} \quad (1459)$$

We see no compelling reason to prefer one particular auxiliary  $V$ -matrix over another. Each has its own coherent interpretations; *e.g.* §5.4.2, §6.6, §B.4.5. Neither can we say that any particular problem formulation produces generally better results than another.<sup>7.6</sup>

## 7.1 First prevalent problem: Projection on PSD cone

This first problem

$$\left. \begin{aligned} & \underset{D}{\text{minimize}} \quad \| -V_N^T (D - H) V_N \|^2_F \\ & \text{subject to} \quad \begin{aligned} & \text{rank } V_N^T D V_N \leq \rho \\ & D \in \mathbb{EDM}^N \end{aligned} \end{aligned} \right\} \quad \text{Problem 1} \quad (1460)$$

poses Euclidean projection of  $-V_N^T H V_N$  (in subspace  $\mathbb{S}^{N-1}$ ) on a generally nonconvex subset (when  $\rho < N-1$ ) of a positive semidefinite cone boundary  $\partial \mathbb{S}_+^{N-1}$  whose elemental matrices have rank no greater than desired affine dimension  $\rho$  (§5.7.1.1). Problem 1 finds the closest EDM  $D$  in the sense of Schoenberg. (1052) [349] As it is stated, this optimization problem is convex only when desired affine dimension is largest  $\rho = N-1$  although its analytical solution is known [287, thm.14.4.2] for all nonnegative  $\rho \leq N-1$ .<sup>7.7</sup>

We assume only that the given measurement matrix  $H$  is symmetric;<sup>7.8</sup>

$$H \in \mathbb{S}^N \quad (1461)$$

Arranging the eigenvalues  $\lambda_i$  of  $-V_N^T H V_N$  in nonincreasing order for all  $i$  ( $\lambda_i \geq \lambda_{i+1}$  with corresponding  $i^{\text{th}}$  eigenvector  $v_i$ ), then an optimal solution to Problem 1 is [394, §2]

$$-V_N^T D^* V_N = \sum_{i=1}^{\rho} \max\{0, \lambda_i\} v_i v_i^T \quad (1462)$$

<sup>7.5</sup>Isomorphism  $T(Y) = V_N^{\dagger T} Y V_N^\dagger$  onto  $\mathbb{S}_c^N = \{V X V^T \mid X \in \mathbb{S}^N\}$  relates the map in (1457) to that in (1458), but is not an isometry. This behavior may be observed via MATLAB program `isedm()` (provided on *Wikimization* [431]) that solves (1456.1) for any desired upper bound on affine dimension  $\rho$  and allows selection of auxiliary matrix  $V$  or  $V_N$ .

<sup>7.6</sup>All four problem formulations (1456) produce identical results when affine dimension  $r$ , implicit to a realizable measurement matrix  $H$ , does not exceed desired affine dimension  $\rho$ ; because, the optimal objective value will vanish ( $\|\star\| = 0$ ).

<sup>7.7</sup>being first pronounced in the context of multidimensional scaling by Mardia [286] in 1978 who attributes the generic result (§7.1.2) to Eckart & Young, 1936 [147].

<sup>7.8</sup>Projection, in Problem 1, is on a rank  $\rho$  subset of positive semidefinite cone  $\mathbb{S}_+^{N-1}$  (§2.9.2.1) in the subspace of symmetric matrices  $\mathbb{S}^{N-1}$ . It is wrong here to zero the main diagonal of given  $H$  because first projecting  $H$ , on the symmetric hollow subspace, places an elbow in the path of projection in Problem 1. (Figure 178)

where

$$-V_N^T H V_N \triangleq \sum_{i=1}^{N-1} \lambda_i v_i v_i^T \in \mathbb{S}^{N-1} \quad (1463)$$

is an eigenvalue decomposition and

$$D^* \in \mathbb{EDM}^N \quad (1464)$$

is an optimal Euclidean distance matrix.

In §7.1.4 we show how to transform Problem 1 to a convex optimization problem for any  $\rho$ .

### 7.1.1 Closest-EDM Problem 1, convex case

**7.1.1.0.1 Proof.** *Solution (1462), convex case.*

When desired affine dimension is unconstrained,  $\rho = N - 1$ , the rank function disappears from (1460) leaving a convex optimization problem; a simple unique minimum-distance projection on positive semidefinite cone  $\mathbb{S}_+^{N-1}$ : *videlicet*

$$\begin{aligned} & \underset{D \in \mathbb{S}_h^N}{\text{minimize}} \quad \| -V_N^T (D - H) V_N \|_F^2 \\ & \text{subject to} \quad -V_N^T D V_N \succeq 0 \end{aligned} \quad (1465)$$

by (1052). Because

$$\mathbb{S}^{N-1} = -V_N^T \mathbb{S}_h^N V_N \quad (1156)$$

then the necessary and sufficient conditions for projection in isometrically isomorphic  $\mathbb{R}^{N(N-1)/2}$  on selfdual (383) positive semidefinite cone  $\mathbb{S}_+^{N-1}$  are:<sup>7.9</sup> (§E.9.2.0.1) (1748) (*confer*(2232))

$$\begin{aligned} & -V_N^T D^* V_N \succeq 0 \\ & -V_N^T D^* V_N (-V_N^T D^* V_N + V_N^T H V_N) = \mathbf{0} \\ & -V_N^T D^* V_N + V_N^T H V_N \succeq 0 \end{aligned} \quad (1466)$$

Symmetric  $-V_N^T H V_N$  is diagonalizable hence decomposable in terms of its eigenvectors  $v$  and eigenvalues  $\lambda$  as in (1463). Therefore (*confer*(1462))

$$-V_N^T D^* V_N = \sum_{i=1}^{N-1} \max\{0, \lambda_i\} v_i v_i^T \quad (1467)$$

satisfies (1466), optimally solving (1465). To see that, recall: these eigenvectors constitute an orthogonal set and

$$-V_N^T D^* V_N + V_N^T H V_N = - \sum_{i=1}^{N-1} \min\{0, \lambda_i\} v_i v_i^T \quad (1468)$$

◆

---

<sup>7.9</sup>The Karush-Kuhn-Tucker (KKT) optimality conditions, [312, p.328] [65, §5.5.3] for problem (1465), are identical to these conditions for projection on a convex cone.

### 7.1.2 generic problem, projection on PSD cone

Prior to determination of  $D^*$ , analytical solution (1462) to Problem 1 is equivalent to solution of a generic rank-constrained projection problem: Given desired affine dimension  $\rho$  and

$$A \triangleq -V_N^T H V_N = \sum_{i=1}^{N-1} \lambda_i v_i v_i^T \in \mathbb{S}^{N-1} \quad (1463)$$

Euclidean projection on a rank  $\rho$  subset of a positive semidefinite cone (on a generally nonconvex subset of the PSD cone boundary  $\partial \mathbb{S}_+^{N-1}$  when  $\rho < N-1$ )

$$\left. \begin{array}{ll} \text{minimize}_{B \in \mathbb{S}^{N-1}} & \|B - A\|_F^2 \\ \text{subject to} & \begin{array}{l} \text{rank } B \leq \rho \\ B \succeq 0 \end{array} \end{array} \right\} \text{ Generic 1 } \quad (1469)$$

has well known optimal solution (Eckart & Young) [147]

$$B^* \triangleq -V_N^T D^* V_N = \sum_{i=1}^{\rho} \max\{0, \lambda_i\} v_i v_i^T \in \mathbb{S}^{N-1} \quad (1462)$$

Once optimal  $B^*$  is found, the technique of §5.12 can be used to determine a uniquely corresponding optimal Euclidean distance matrix  $D^*$ ; a unique correspondence by injectivity arguments in §5.6.2.

#### 7.1.2.1 Projection on rank $\rho$ subset of PSD cone

Because (1156) provides invertible mapping to the generic problem, then Problem 1

$$\left. \begin{array}{ll} \text{minimize}_{D \in \mathbb{S}_h^N} & \| -V_N^T (D - H) V_N \|_F^2 \\ \text{subject to} & \begin{array}{l} \text{rank } V_N^T D V_N \leq \rho \\ -V_N^T D V_N \succeq 0 \end{array} \end{array} \right\} \quad (1470)$$

is truly a Euclidean projection of vectorized  $-V_N^T H V_N$  on that generally nonconvex subset of symmetric matrices (belonging to positive semidefinite cone  $\mathbb{S}_+^{N-1}$ ) having rank no greater than desired affine dimension  $\rho$ ; <sup>7.10</sup> called *rank  $\rho$  subset*: (265)

$$\mathbb{S}_+^{N-1} \setminus \mathbb{S}_+^{N-1}(\rho + 1) = \{X \in \mathbb{S}_+^{N-1} \mid \text{rank } X \leq \rho\} \quad (220)$$

### 7.1.3 Choice of spectral cone

Spectral projection substitutes projection on a polyhedral cone, containing a complete set of eigenspectra (§5.11.1.0.3), in place of projection on a convex set of diagonalizable matrices; *e.g.*, (1483). In this section we develop a method of spectral projection for constraining rank of positive semidefinite matrices in a proximity problem like (1469). We will see why an orthant turns out to be the best choice of spectral cone, and why presorting is critical.

Define a nonlinear permutation-operator

$$\pi(x) : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad (1471)$$

that sorts its vector argument  $x$  into nonincreasing order; *a.k.a., presorting function*.

---

<sup>7.10</sup>Recall: affine dimension is a lower bound on embedding (§2.3.1), equal to dimension of the smallest affine set in which points from a list  $X$  corresponding to an EDM  $D$  can be embedded.

### 7.1.3.0.1 Definition. Spectral projection.

Let  $R$  be an orthogonal matrix and  $\Lambda$  a nonincreasingly ordered diagonal matrix of eigenvalues. *Spectral projection* means unique minimum-distance projection of a rotated  $(R, \S B.5.5)$  nonincreasingly ordered  $(\pi)$  vector  $(\delta)$  of eigenvalues

$$\pi(\delta(R^T \Lambda R)) \quad (1472)$$

on a polyhedral cone containing all eigenspectra corresponding to a rank  $\rho$  subset of a PSD cone ([§2.9.2.1](#)) or the EDM cone (in Cayley-Menger form, [§5.11.2.3](#)).  $\triangle$

In the simplest and most common case, projection on a positive semidefinite cone, orthogonal matrix  $R$  equals  $I$  ([§7.1.4.0.1](#)) and diagonal matrix  $\Lambda$  is ordered during diagonalization ([§A.5.1](#)). Then spectral projection simply means projection of  $\delta(\Lambda)$  on a subset of the nonnegative orthant, as we shall now ascertain:

It is curious how nonconvex Problem 1 has such a simple analytical solution ([\(1462\)](#)). Although solution to generic problem ([\(1469\)](#)) is well known since 1936 [[147](#)], its equivalence was observed in 1997 [[394](#), §2] to projection of an ordered vector of eigenvalues (in diagonal matrix  $\Lambda$ ) on a subset of the monotone nonnegative cone ([§2.13.10.4.2](#))

$$\mathcal{K}_{\mathcal{M}+} = \{v \mid v_1 \geq v_2 \geq \cdots \geq v_{N-1} \geq 0\} \subseteq \mathbb{R}_+^{N-1} \quad (434)$$

Of interest, momentarily, is only the smallest convex subset of the monotone nonnegative cone  $\mathcal{K}_{\mathcal{M}+}$  containing every nonincreasingly ordered eigenspectrum corresponding to a rank  $\rho$  subset of positive semidefinite cone  $\mathbb{S}_+^{N-1}$ ; *id est*,

$$\mathcal{K}_{\mathcal{M}+}^\rho \triangleq \{v \in \mathbb{R}^\rho \mid v_1 \geq v_2 \geq \cdots \geq v_\rho \geq 0\} \subseteq \mathbb{R}_+^\rho \quad (1473)$$

a pointed polyhedral cone, a  $\rho$ -dimensional convex subset of the monotone nonnegative cone  $\mathcal{K}_{\mathcal{M}+} \subseteq \mathbb{R}_+^{N-1}$  having property, for  $\lambda$  denoting eigenspectra,

$$\begin{bmatrix} \mathcal{K}_{\mathcal{M}+}^\rho \\ \mathbf{0} \end{bmatrix} = \pi(\lambda(\text{rank } \rho \text{ subset})) \subseteq \mathcal{K}_{\mathcal{M}+}^{N-1} \triangleq \mathcal{K}_{\mathcal{M}+} \quad (1474)$$

For each and every elemental eigenspectrum

$$\gamma \in \lambda(\text{rank } \rho \text{ subset}) \subseteq \mathbb{R}_+^{N-1} \quad (1475)$$

of the rank  $\rho$  subset (ordered or unordered in  $\lambda$ ), there is a nonlinear surjection  $\pi(\gamma)$  onto  $\mathcal{K}_{\mathcal{M}+}^\rho$ .

### 7.1.3.0.2 Exercise. Smallest spectral cone.

Prove that there is no convex subset of  $\mathcal{K}_{\mathcal{M}+}$  smaller than  $\mathcal{K}_{\mathcal{M}+}^\rho$  containing every ordered eigenspectrum corresponding to the rank  $\rho$  subset of a positive semidefinite cone ([§2.9.2.1](#)).  $\blacktriangledown$

**7.1.3.0.3 Proposition.** (Hardy-Littlewood-Pólya) *Inequalities.* [203, §X] [58, §1.2]  
Any vectors  $\sigma$  and  $\gamma$  in  $\mathbb{R}^{N-1}$  satisfy a tight inequality

$$\pi(\sigma)^T \pi(\gamma) \geq \sigma^T \gamma \geq \pi(\sigma)^T \Xi \pi(\gamma) \quad (1476)$$

where  $\Xi$  is the order-reversing permutation matrix defined in ([1900](#)), and permutator  $\pi(\gamma)$  is a nonlinear function that sorts vector  $\gamma$  into nonincreasing order thereby providing the greatest upper bound and least lower bound with respect to every possible sorting.  $\diamond$

**7.1.3.0.4 Corollary.** *Monotone nonnegative sort.*

Any given vectors  $\sigma, \gamma \in \mathbb{R}^{N-1}$  satisfy a tight Euclidean distance inequality

$$\|\pi(\sigma) - \pi(\gamma)\| \leq \|\sigma - \gamma\| \quad (1477)$$

where nonlinear function  $\pi(\gamma)$  sorts vector  $\gamma$  into nonincreasing order thereby providing the least lower bound with respect to every possible sorting.  $\diamond$

Given  $\gamma \in \mathbb{R}^{N-1}$

$$\inf_{\sigma \in \mathbb{R}_+^{N-1}} \|\sigma - \gamma\| = \inf_{\sigma \in \mathbb{R}_+^{N-1}} \|\pi(\sigma) - \pi(\gamma)\| = \inf_{\sigma \in \mathbb{R}_+^{N-1}} \|\sigma - \pi(\gamma)\| = \inf_{\sigma \in \mathcal{K}_{\mathcal{M}+}} \|\sigma - \pi(\gamma)\| \quad (1478)$$

Yet for  $\gamma$  representing an arbitrary vector of eigenvalues, because

$$\inf_{\substack{\sigma \in [\mathbb{R}_+^\rho] \\ \mathbf{0}}} \|\sigma - \gamma\|^2 \geq \inf_{\substack{\sigma \in [\mathbb{R}_+^\rho] \\ \mathbf{0}}} \|\sigma - \pi(\gamma)\|^2 = \inf_{\substack{\sigma \in [\mathcal{K}_{\mathcal{M}+}^\rho] \\ \mathbf{0}}} \|\sigma - \pi(\gamma)\|^2 \quad (1479)$$

then projection of  $\gamma$  on the eigenspectra of a rank  $\rho$  subset can be tightened simply by presorting  $\gamma$  into nonincreasing order.

**Proof.** Simply because  $\pi(\gamma)_{1:\rho} \succeq \pi(\gamma_{1:\rho})$

$$\begin{aligned} \inf_{\substack{\sigma \in [\mathbb{R}_+^\rho] \\ \mathbf{0}}} \|\sigma - \gamma\|^2 &= \gamma_{\rho+1:N-1}^T \gamma_{\rho+1:N-1} + \inf_{\sigma \in \mathbb{R}_+^{N-1}} \|\sigma_{1:\rho} - \gamma_{1:\rho}\|^2 \\ &= \gamma^T \gamma + \inf_{\sigma \in \mathbb{R}_+^{N-1}} \sigma_{1:\rho}^T \sigma_{1:\rho} - 2\sigma_{1:\rho}^T \gamma_{1:\rho} \\ &\geq \gamma^T \gamma + \inf_{\sigma \in \mathbb{R}_+^{N-1}} \sigma_{1:\rho}^T \sigma_{1:\rho} - 2\sigma_{1:\rho}^T \pi(\gamma)_{1:\rho} \end{aligned} \quad (1480)$$

$$\inf_{\substack{\sigma \in [\mathbb{R}_+^\rho] \\ \mathbf{0}}} \|\sigma - \gamma\|^2 \geq \inf_{\substack{\sigma \in [\mathbb{R}_+^\rho] \\ \mathbf{0}}} \|\sigma - \pi(\gamma)\|^2$$

$\blacklozenge$

**7.1.3.1 Orthant is best spectral cone for Problem 1**

This means unique minimum-distance projection of  $\gamma$  on the nearest spectral member of the rank  $\rho$  subset is tantamount to presorting  $\gamma$  into nonincreasing order. Only then does unique spectral projection on a subset  $\mathcal{K}_{\mathcal{M}+}^\rho$  of the monotone nonnegative cone become equivalent to unique spectral projection on a subset  $\mathbb{R}_+^\rho$  of the nonnegative orthant (which is simpler); in other words, unique minimum-distance projection of sorted  $\gamma$  on the nonnegative orthant in a  $\rho$ -dimensional subspace of  $\mathbb{R}^N$  is indistinguishable from its projection on the subset  $\mathcal{K}_{\mathcal{M}+}^\rho$  of the monotone nonnegative cone in that same subspace.

**7.1.4 Closest-EDM Problem 1, “nonconvex” case**

Proof of solution (1462), for projection on a rank  $\rho$  subset of positive semidefinite cone  $\mathbb{S}_+^{N-1}$ , can be algebraic in nature. [394, §2] Here we derive that known result but instead using a more geometric argument via spectral projection on a polyhedral cone (subsuming the proof in §7.1.1). In so doing, we demonstrate how nonconvex Problem 1 is transformed to a convex optimization:

**7.1.4.0.1 Proof.** *Solution (1462), nonconvex case.*

As explained in §7.1.2, we may instead work with the more facile generic problem (1469). With diagonalization of unknown

$$B \triangleq U\Upsilon U^T \in \mathbb{S}^{N-1} \quad (1481)$$

given desired affine dimension  $0 \leq \rho \leq N-1$  and diagonalizable

$$A \triangleq Q\Lambda Q^T \in \mathbb{S}^{N-1} \quad (1482)$$

having eigenvalues in  $\Lambda$  arranged in nonincreasing order, by (49) the generic problem is equivalent to

$$\begin{array}{lll} \underset{\substack{B \in \mathbb{S}^{N-1} \\ \text{subject to} \\ B \succeq 0}}{\text{minimize}} & \|B - A\|_F^2 & \underset{\substack{R, \Upsilon \\ \text{subject to} \\ \Upsilon \succeq 0 \\ R^{-1} = R^T}}{\text{minimize}} & \|\Upsilon - R^T \Lambda R\|_F^2 \\ & \text{rank } B \leq \rho & & \text{rank } \Upsilon \leq \rho \end{array} \quad (1483)$$

where

$$R \triangleq Q^T U \in \mathbb{R}^{N-1 \times N-1} \quad (1484)$$

is a bijection in  $U$  on the set of orthogonal matrices. We propose solving (1483) by instead solving the problem sequence:

$$\begin{array}{ll} \underset{\substack{\Upsilon \\ \text{subject to} \\ \Upsilon \succeq 0}}{\text{minimize}} & \|\Upsilon - R^T \Lambda R\|_F^2 \\ & \text{rank } \Upsilon \leq \rho & \text{(a)} \\ \underset{\substack{R \\ \text{subject to} \\ R^{-1} = R^T}}{\text{minimize}} & \|\Upsilon^* - R^T \Lambda R\|_F^2 \\ & & \text{(b)} \end{array} \quad (1485)$$

Problem (1485a) is equivalent to:

- (1) orthogonal projection of  $R^T \Lambda R$  on an  $N-1$ -dimensional subspace of isometrically isomorphic  $\mathbb{R}^{N(N-1)/2}$  containing  $\delta(\Upsilon) \in \mathbb{R}_+^{N-1}$
- (2) nonincreasingly ordering the result,
- (3) unique minimum-distance projection of the ordered result on  $\begin{bmatrix} \mathbb{R}_+^\rho \\ \mathbf{0} \end{bmatrix}$  (§E.9.5).

Projection on that  $N-1$ -dimensional subspace amounts to zeroing  $R^T \Lambda R$  at all entries off the main diagonal; thus, the equivalent sequence leading with a spectral projection:

$$\begin{array}{ll} \underset{\substack{\Upsilon \\ \text{subject to} \\ \delta(\Upsilon) \in \begin{bmatrix} \mathbb{R}_+^\rho \\ \mathbf{0} \end{bmatrix}}}{\text{minimize}} & \|\delta(\Upsilon) - \pi(\delta(R^T \Lambda R))\|^2 \\ & \text{(a)} \\ \underset{\substack{R \\ \text{subject to} \\ R^{-1} = R^T}}{\text{minimize}} & \|\Upsilon^* - R^T \Lambda R\|_F^2 \\ & \text{(b)} \end{array} \quad (1486)$$

Because any permutation matrix is an orthogonal matrix,  $\delta(R^T \Lambda R) \in \mathbb{R}^{N-1}$  can always be arranged in nonincreasing order without loss of generality; hence, permutation operator  $\pi$ .

Unique minimum-distance projection of vector  $\pi(\delta(R^T \Lambda R))$  on the  $\rho$ -dimensional subset  $\begin{bmatrix} \mathbb{R}_+^\rho \\ \mathbf{0} \end{bmatrix}$  of nonnegative orthant  $\mathbb{R}_+^{N-1}$  requires: (§E.9.2.0.1)

$$\begin{aligned} \delta(\Upsilon^*)_{\rho+1:N-1} &= \mathbf{0} \\ \delta(\Upsilon^*) &\succeq 0 \\ \delta(\Upsilon^*)^T (\delta(\Upsilon^*) - \pi(\delta(R^T \Lambda R))) &= 0 \\ \delta(\Upsilon^*) - \pi(\delta(R^T \Lambda R)) &\succeq 0 \end{aligned} \tag{1487}$$

which are necessary and sufficient conditions. Any value  $\Upsilon^*$  satisfying conditions (1487) is optimal for (1486a). So

$$\delta(\Upsilon^*)_i = \begin{cases} \max \left\{ 0, \pi(\delta(R^T \Lambda R))_i \right\}, & i=1 \dots \rho \\ 0, & i=\rho+1 \dots N-1 \end{cases} \tag{1488}$$

specifies an optimal solution. The lower bound on the objective with respect to  $R$  in (1486b) is tight: by (1453)

$$\| |\Upsilon^*| - |\Lambda| \|_F \leq \| \Upsilon^* - R^T \Lambda R \|_F \tag{1489}$$

where  $| |$  denotes absolute entry-value. For selection of  $\Upsilon^*$  as in (1488), this lower bound is attained when (*confer* §C.4.2.2)

$$R^* = I \tag{1490}$$

which is the known solution. ♦

#### 7.1.4.1 significance

Importance of this well-known [147] optimal solution (1462) for projection on a rank  $\rho$  subset of a positive semidefinite cone should not be dismissed:

- Problem 1 (1460) and its generic form (1469), as stated, are generally nonconvex. Their known analytical solution encompasses projection on a rank  $\rho$  subset (220) of a positive semidefinite cone (generally, a nonconvex subset of its boundary) from either the exterior or interior of that cone.<sup>7.11</sup> By problem transformation to the spectral domain, projection on a rank  $\rho$  subset becomes a convex optimization problem.
- This solution is closed form.
- This solution is equivalent to projection on a polyhedral cone in the spectral domain (spectral projection §7.1.3.0.1, projection on a spectral cone §5.11.1.0.2); a necessary and sufficient condition (§A.3.1) for membership of a symmetric matrix to a rank  $\rho$  subset of a positive semidefinite cone (§2.9.2.1).
- A minimum-distance projection, on a rank  $\rho$  subset of a positive semidefinite cone, is a positive semidefinite matrix orthogonal (in the Euclidean sense) to direction of projection<sup>7.12</sup> because  $U^* = Q$  in (1484).
- For the convex case problem (1465), this solution is always unique. Otherwise, distinct eigenvalues (multiplicity 1) in  $\Lambda$  guarantee uniqueness of this solution by the reasoning in §A.5.0.1.<sup>7.13</sup>

<sup>7.11</sup>Projection on the boundary from the interior, of a convex Euclidean body, is generally a nonconvex problem. (§E.9.1.1.2)

<sup>7.12</sup>But Theorem E.9.2.0.1, for unique projection on a closed convex cone, does not apply here because direction of projection is not necessarily a member of the dual PSD cone. This occurs, for example, whenever positive eigenvalues are truncated.

<sup>7.13</sup>Uncertainty of uniqueness prevents the erroneous conclusion that a rank  $\rho$  subset (220) were a convex body by the *Bunt-Motzkin theorem* (§E.9.0.0.1).

### 7.1.4.2 list projection interpretation

Because  $-VDV\frac{1}{2} = X^T X$  when point list  $X$  is geometrically centered,  $X\mathbf{1} = \mathbf{0}$ , Problem 1 can be equivalently restated: by (1054)

$$(1) \quad \begin{array}{lll} \underset{D}{\text{minimize}} & \| -V(D-H)V \|^2_F \\ \text{subject to} & \text{rank } VDV \leq \rho & \equiv \\ & D \in \mathbb{EDM}^N & \min_{X \in \mathbb{R}^{N \times N}} \| X^T X - Y^T Y \|^2_F \quad (\text{G}) \end{array} \quad (1491)$$

where  $Y \in \mathbb{R}^{N \times N}$  comprises geometrically centered point list estimates ( $Y = YV$ ) whose dimensionality is to be reduced (by best fit) to

$$\rho \leq \eta \triangleq \min\{n, N\} \quad (1492)$$

We call (1491.G) the *Gram-form* Problem 1; it may be interpreted as minimum-distance projection of  $Y^T Y \in \mathbb{S}^N$  on a rank  $\rho$  subset (§2.9.2.1) of the PSD cone in isometrically isomorphic  $\mathbb{R}^{N(N+1)/2}$ . Geometrically centered  $Y$  remains centered, postprojection, because the subspace  $\mathbb{S}_c^N$  of symmetric geometrically centered matrices  $VY^T YV$  (1135) is invariant to projection on a positive semidefinite cone by Lemma 6.8.1.1.1.

Orthogonal projection of estimates  $Y$ , on span of  $\rho$  principal eigenvectors of  $YY^T \in \mathbb{S}^n$ , provides unique (not rotation invariant) optimal  $X^*$  in the sense

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{n \times N}}{\text{minimize}} & \| X^T X - Y^T Y \|^2_F = \underset{X \in \mathbb{R}^{n \times N}}{\text{minimize}} \| XX^T - YY^T \|^2_F = \sum_{i=\rho+1}^n \lambda(Y^T Y)_i^2 \\ \text{subject to} & \text{rank}(X^T X) \leq \rho \quad \text{subject to} \quad \text{rank}(XX^T) \leq \rho \end{array} \quad (1493)$$

where  $\text{rank}(X^T X) = \text{rank}(XX^T)$  (1618). Defining nonincreasingly ordered diagonalization  $YY^T \triangleq Q_n \Lambda_n Q_n^T \in \mathbb{S}^n$ , then orthogonal projection of  $Y$  is (§E.3.2) 7.14

$$X^* = Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T Y \in \mathbb{R}^{n \times N} \quad (1494)$$

So

$$\begin{aligned} \| X^* X^{*\top} - YY^T \|^2_F &= \| Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T YY^T - YY^T \|^2_F \\ &= \left\| \begin{bmatrix} \Lambda_n(1:\rho, 1:\rho) & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} \end{bmatrix} - \Lambda_n \right\|_F^2 = \sum_{i=\rho+1}^n \lambda(Y^T Y)_i^2 \end{aligned} \quad (1495)$$

projection of list  $Y$  on a subspace solves projection of Gram matrix  $Y^T Y$  [sic] on a positive semidefinite cone (1493); quite a remarkable interpretation. 7.15 [287, §14.4] [328, §2]

Yet there is a more plain interpretation:

$$X^* X^{*\top} = Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T YY^T = Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T YY^T Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T \in \mathbb{S}_+^n \quad (1496)$$

is the orthogonal projection of  $YY^T$  on the closest  $\rho(\rho+1)/2$ -dimensional subspace

$$Q_n(:, 1:\rho) \mathbb{S}^\rho Q_n(:, 1:\rho)^T = Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T \mathbb{S}^n Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T \quad (1497)$$

of a rotated Cartesian coordinate system  $Q_n \mathbb{S}^n Q_n^T$  in isomorphic  $\mathbb{R}^{n(n+1)/2}$ . That it is the closest subspace, comes from §2.13.8.1.1. That (1497) is the smallest subspace containing the smallest face (that contains  $X^* X^{*\top}$ ) of PSD cone  $\mathbb{S}_+^n$ , is a result from §2.9.2.4.

7.14 Reconstruction of  $X^*$ , with dimension  $\rho$  instead of  $n$ , is disclosed in §5.12.2.

7.15 This might imply existence of an isomorphism (§2.2.1.0.1) relating vector space  $\mathbb{R}^{nN}$  (containing vectorized list  $X$ ) to vector space  $\mathbb{R}^{N(N+1)/2}$  (containing vectorized cone  $\mathbb{S}_+^N$ ); but there is none. Such an isomorphism might be an isometry (2.2.1.1.1) were  $\| X^* X^{*\top} - YY^T \|^2_F$  equal to

$$\begin{array}{lll} \underset{X \in \mathbb{R}^{n \times N}}{\text{minimize}} & \| X - Y \|^2_F &= \| Q_n(:, 1:\rho) Q_n(:, 1:\rho)^T Y - Y \|^2_F &= \sum_{i=\rho+1}^n \lambda(Y^T Y)_i \\ \text{subject to} & X\mathbf{1} = \mathbf{0} & & \\ & \text{rank } X \leq \rho & & \end{array}$$

but the square of eigenvalues is absent with respect to (1495); numerically verifiable by means of problem transformation in §4.9 and a few convex iterations (§4.5.1).

### 7.1.5 Problem 1 in spectral norm, convex case

When instead we pose the matrix 2-norm (spectral norm) in Problem 1 (1460) for the convex case  $\rho = N - 1$ , then the new problem

$$\begin{aligned} & \underset{D}{\text{minimize}} && \| -V_N^T(D - H)V_N \|_2 \\ & \text{subject to} && D \in \mathbb{EDM}^N \end{aligned} \quad (1498)$$

is convex although its solution is not necessarily unique;<sup>7.16</sup> giving rise to nonorthogonal projection (§E.1) on positive semidefinite cone  $\mathbb{S}_+^{N-1}$ . Indeed, its solution set includes the Frobenius solution (1462) for the convex case whenever  $-V_N^T H V_N$  is a normal matrix. [210, §1] [201] [65, §8.1.1] Proximity problem (1498) is equivalent to

$$\begin{aligned} & \underset{\mu, D}{\text{minimize}} && \mu \\ & \text{subject to} && -\mu I \preceq -V_N^T(D - H)V_N \preceq \mu I \\ & && D \in \mathbb{EDM}^N \end{aligned} \quad (1499)$$

by (1876) where

$$\mu^* = \max_i \{ |\lambda(-V_N^T(D^* - H)V_N)_i|, i = 1 \dots N-1 \} \in \mathbb{R}_+ \quad (1500)$$

is the minimized largest absolute eigenvalue (due to matrix symmetry).

For lack of unique solution here, we prefer the Frobenius rather than spectral norm.

## 7.2 Second prevalent problem: Projection on EDM cone in $\sqrt{d_{ij}}$

Let

$$\sqrt{D} \triangleq [\sqrt{d_{ij}}] \in \mathcal{K} = \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N} \quad (1501)$$

be an unknown matrix of absolute distance; *id est*,

$$D = [d_{ij}] \triangleq \sqrt[3]{D} \circ \sqrt[3]{D} \in \mathbb{EDM}^N \quad (1502)$$

where  $\circ$  denotes Hadamard product. The second prevalent proximity problem is a Euclidean projection (in the natural coordinates  $\sqrt{d_{ij}}$ ) of matrix  $H$  on a nonconvex subset of the boundary of the nonconvex cone of Euclidean absolute-distance matrices  $\text{rel } \partial \sqrt{\mathbb{EDM}^N}$ : (§6.3, confer Figure 163b)

$$\left. \begin{aligned} & \underset{\sqrt{D}}{\text{minimize}} && \|\sqrt[3]{D} - H\|_F^2 \\ & \text{subject to} && \text{rank } V_N^T D V_N \leq \rho \\ & && \sqrt{D} \in \sqrt{\mathbb{EDM}^N} \end{aligned} \right\} \quad \text{Problem 2} \quad (1503)$$

where

$$\sqrt{\mathbb{EDM}^N} = \{ \sqrt[3]{D} \mid D \in \mathbb{EDM}^N \} \quad (1328)$$

This statement of the second proximity problem is considered difficult to solve because of the constraint on desired affine dimension  $\rho$  (§5.7.2) and because the objective function

$$\|\sqrt[3]{D} - H\|_F^2 = \sum_{i,j} (\sqrt{d_{ij}} - h_{ij})^2 \quad (1504)$$

---

<sup>7.16</sup>For each and every  $|t| \leq 2$ , for example,  $\begin{bmatrix} 2 & 0 \\ 0 & t \end{bmatrix}$  has the same spectral-norm value.

is expressed in the natural coordinates; projection on a doubly nonconvex set.

Our solution to this second problem prevalent in the literature requires measurement matrix  $H$  to be nonnegative;

$$H = [h_{ij}] \in \mathbb{R}_+^{N \times N} \quad (1505)$$

If the  $H$  matrix given has negative entries, then the technique of solution presented here becomes invalid. As explained in §7.0.1, projection of  $H$  on  $\mathcal{K} = \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N}$  (1447) prior to application of this proposed solution is incorrect.

### 7.2.1 Convex case

When  $\rho = N - 1$ , the rank constraint vanishes and a convex problem that is equivalent to (1454) emerges:<sup>7.17</sup>

$$\begin{array}{ll} \underset{\sqrt{D}}{\text{minimize}} & \|\sqrt{D} - H\|_F^2 \\ \text{subject to} & \sqrt{D} \in \sqrt{\text{EDM}^N} \end{array} \Leftrightarrow \begin{array}{ll} \underset{D}{\text{minimize}} & \sum_{i,j} d_{ij} - 2h_{ij}\sqrt{d_{ij}} + h_{ij}^2 \\ \text{subject to} & D \in \text{EDM}^N \end{array} \quad (1506)$$

For any fixed  $i$  and  $j$ , the argument of summation is a convex function of  $d_{ij}$  because (for nonnegative constant  $h_{ij}$ ) the negative square root is convex in nonnegative  $d_{ij}$  and because  $d_{ij} + h_{ij}^2$  is affine (convex). Because the sum of any number of convex functions in  $D$  remains convex [65, §3.2.1] and because the feasible set is convex in  $D$ , we have a convex optimization problem:

$$\begin{array}{ll} \underset{D}{\text{minimize}} & \mathbf{1}^T(D - 2H \circ \sqrt{D})\mathbf{1} + \|H\|_F^2 \\ \text{subject to} & D \in \text{EDM}^N \end{array} \quad (1507)$$

The objective function being a sum of strictly convex functions is, moreover, strictly convex in  $D$  on the nonnegative orthant. Existence of a unique solution  $D^*$  for this second prevalent problem depends upon nonnegativity of  $H$  and a convex feasible set (§3.1.1).<sup>7.18</sup>

#### 7.2.1.1 Equivalent semidefinite program, Problem 2, convex case

Convex problem (1506) is numerically solvable for its global minimum using an interior-point method [453] [321] [309] [444] [12] [165]. We translate (1506) to an equivalent semidefinite program (SDP) for a pedagogical reason made clear in §7.2.2.2 and because there exist readily available computer programs for numerical solution [191] [446] [447] [400] [35] [445] [389] [373].

Substituting a new matrix variable  $Y \triangleq [y_{ij}] \in \mathbb{R}_+^{N \times N}$

$$h_{ij}\sqrt{d_{ij}} \leftarrow y_{ij} \quad (1508)$$

Boyd proposes: problem (1506) is equivalent to the semidefinite program

$$\begin{array}{ll} \underset{D, Y}{\text{minimize}} & \sum_{i,j} d_{ij} - 2y_{ij} + h_{ij}^2 \\ \text{subject to} & \begin{bmatrix} d_{ij} & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0, \quad i, j = 1 \dots N \\ & D \in \text{EDM}^N \end{array} \quad (1509)$$

<sup>7.17</sup> still thought to be a nonconvex problem as late as 1997 [395] even though discovered convex by de Leeuw in 1993. [115] [56, §13.6] Yet using methods from §3, it can be easily ascertained:  $\|\sqrt{D} - H\|_F$  is not convex in  $D$ .

<sup>7.18</sup>The transformed problem in variable  $D$  no longer describes Euclidean projection on an EDM cone. Otherwise we might erroneously conclude  $\sqrt{\text{EDM}^N}$  were a convex body by the *Bunt-Motzkin theorem* (§E.9.0.0.1).

To see that, recall:  $d_{ij} \geq 0$  is implicit to  $D \in \mathbb{EDM}^N$  (§5.8.1, (1052)). So when  $H \in \mathbb{R}_+^{N \times N}$  is nonnegative, as assumed,

$$\begin{bmatrix} d_{ij} & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0 \Leftrightarrow h_{ij} \sqrt{d_{ij}} \geq \sqrt{y_{ij}^2} \quad (1510)$$

Minimization of the objective function implies maximization of  $y_{ij}$  that is bounded above. Hence nonnegativity of  $y_{ij}$  is implicit to (1509) and, as desired,  $y_{ij} \rightarrow h_{ij} \sqrt{d_{ij}}$  as optimization proceeds. ♦

If the given matrix  $H$  is now assumed symmetric and nonnegative,

$$H = [h_{ij}] \in \mathbb{S}^N \cap \mathbb{R}_+^{N \times N} \quad (1511)$$

then  $Y = H \circ \sqrt[D]{D}$  must belong to  $\mathcal{K} = \mathbb{S}_h^N \cap \mathbb{R}_+^{N \times N}$  (1447). Because  $Y \in \mathbb{S}_h^N$  (§B.4.2 no.20), then

$$\|\sqrt[D]{D} - H\|_F^2 = \sum_{i,j} d_{ij} - 2y_{ij} + h_{ij}^2 = -N \operatorname{tr}(V(D - 2Y)V) + \|H\|_F^2 \quad (1512)$$

So convex problem (1509) is equivalent to the semidefinite program

$$\begin{aligned} & \underset{D, Y}{\text{minimize}} \quad -\operatorname{tr}(V(D - 2Y)V) \\ & \text{subject to} \quad \begin{bmatrix} d_{ij} & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N-1 \\ & \quad Y \in \mathbb{S}_h^N \\ & \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1513)$$

where the constants  $h_{ij}^2$  and  $N$  have been dropped arbitrarily from the objective.

### 7.2.1.2 Gram-form semidefinite program, Problem 2, convex case

There is great advantage to expressing problem statement (1513) in Gram-form because Gram matrix  $G$  is a bidirectional bridge between point list  $X$  and distance matrix  $D$ ; e.g., §5.4.2.2.8, §6.7.0.0.1. This way, problem convexity can be maintained while simultaneously constraining point list  $X$ , Gram matrix  $G$ , and distance matrix  $D$  at our discretion.

Convex problem (1513) may be equivalently written via linear bijective (§5.6.1) EDM operator  $\mathbf{D}(G)$  (1045);

$$\begin{aligned} & \underset{G \in \mathbb{S}_c^N, Y \in \mathbb{S}_h^N}{\text{minimize}} \quad -\operatorname{tr}(V(\mathbf{D}(G) - 2Y)V) \\ & \text{subject to} \quad \begin{bmatrix} \langle \Phi_{ij}, G \rangle & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N-1 \\ & \quad G \succeq 0 \end{aligned} \quad (1514)$$

where distance-square  $D = [d_{ij}] \in \mathbb{S}_h^N$  (1029) is related to  $G = [g_{ij}] \in \mathbb{S}_c^N \cap \mathbb{S}_+^N$  Gram matrix entries by

$$\begin{aligned} d_{ij} &= g_{ii} + g_{jj} - 2g_{ij} \\ &= \langle \Phi_{ij}, G \rangle \end{aligned} \quad (1044)$$

where

$$\Phi_{ij} = (e_i - e_j)(e_i - e_j)^T \in \mathbb{S}_+^N \quad (1031)$$

Confinement of  $G$  to the geometric center subspace provides numerical stability and no loss of generality (*confer*(1392)); implicit constraint  $G\mathbf{1} = \mathbf{0}$  is otherwise unnecessary.

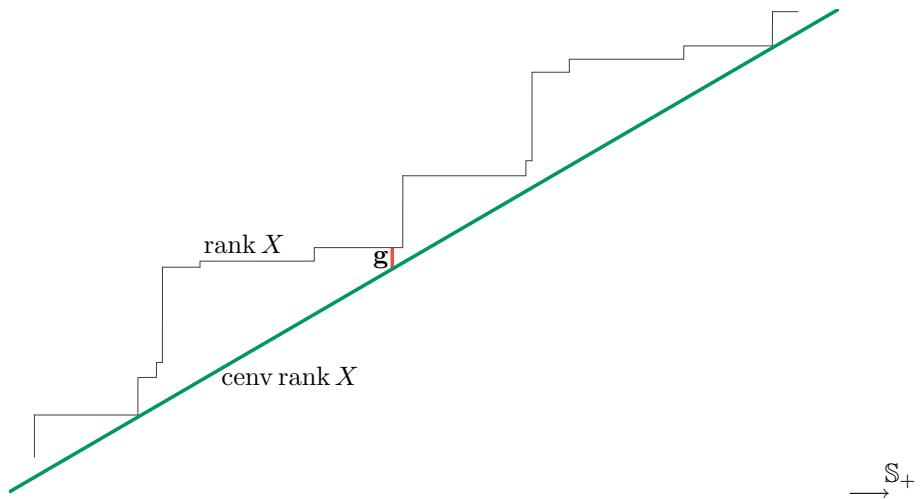


Figure 179: Abstraction of convex envelope of rank function. Rank is a quasiconcave monotonic function on a positive semidefinite cone  $\mathbb{S}_+$ , but its convex envelope is the largest convex function whose epigraph contains it. Vertical bar labelled  $g$  measures a trace–rank gap; *id est*, rank found always exceeds estimate; large decline in trace required here for only a small decrease in rank.

To include constraints on the list  $X \in \mathbb{R}^{n \times N}$ , we would first rewrite (1514)

$$\begin{aligned}
 & \underset{G \in \mathbb{S}_c^N, Y \in \mathbb{S}_h^N, X \in \mathbb{R}^{n \times N}}{\text{minimize}} \quad -\text{tr}(V(\mathbf{D}(G) - 2Y)V) \\
 & \text{subject to} \quad \begin{bmatrix} \langle \Phi_{ij}, G \rangle & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N-1 \\
 & \quad \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} \succeq 0 \\
 & \quad X \in \mathcal{C}
 \end{aligned} \tag{1515}$$

and then introduce the constraints, realized here in abstract membership to some convex set  $\mathcal{C}$ . This problem realization includes a convex relaxation of the nonconvex constraint  $G = X^T X$ . If desired, more constraints on  $G$  could be introduced. These techniques are discussed in §5.4.2.2.8.

## 7.2.2 Minimization of affine dimension in Problem 2

When desired affine dimension  $\rho$  is diminished, the rank function becomes reinserted into problem (1509) that is then rendered difficult to solve because feasible set  $\{D, Y\}$  loses convexity in  $\mathbb{S}_h^N \times \mathbb{R}^{N \times N}$ . Indeed, the rank function is quasiconcave (§3.14) on a positive semidefinite cone; (§2.9.2.9.2) *id est*, its sublevel sets are not convex.

### 7.2.2.1 Rank minimization heuristic

A remedy developed in [293] [153] [154] [152] introduces convex envelope of the quasiconcave rank function: (Figure 179)

**7.2.2.1.1 Definition.** *Convex envelope.*

[224]

Convex envelope  $\text{cenv } f$  of a function  $f: \mathcal{C} \rightarrow \mathbb{R}$  is defined to be the largest convex function  $g$  such that  $g \leq f$  on convex domain  $\mathcal{C} \subseteq \mathbb{R}^n$ . <sup>7.19</sup>  $\triangle$

- [153] [152] Convex envelope of rank function: for  $\sigma_i$  a singular value, (1719)

$$\text{cenv}(\text{rank } A) \text{ on } \{A \in \mathbb{R}^{m \times n} \mid \|A\|_2 \leq \kappa\} = \frac{1}{\kappa} \mathbf{1}^T \sigma(A) = \frac{1}{\kappa} \text{tr} \sqrt{A^T A} \quad (1516)$$

$$\text{cenv}(\text{rank } A) \text{ on } \{A \text{ normal} \mid \|A\|_2 \leq \kappa\} = \frac{1}{\kappa} \|\lambda(A)\|_1 = \frac{1}{\kappa} \text{tr} \sqrt{A^T A} \quad (1517)$$

$$\text{cenv}(\text{rank } A) \text{ on } \{A \in \mathbb{S}_+^n \mid \|A\|_2 \leq \kappa\} = \frac{1}{\kappa} \mathbf{1}^T \lambda(A) = \frac{1}{\kappa} \text{tr}(A) \quad (1518)$$

A properly scaled trace thus represents the best convex lower bound on rank for positive semidefinite matrices. The idea, then, is to substitute convex envelope for rank of some variable  $A \in \mathbb{S}_+^M$  ([§A.6.3.2](#))

$$\text{rank } A \leftarrow \text{cenv}(\text{rank } A) \propto \text{tr } A = \sum_i \sigma(A)_i = \sum_i \lambda(A)_i \quad (1519)$$

which is equivalent to the sum of all eigenvalues or singular values.

- [152] Convex envelope of the cardinality function is proportional to the 1-norm:

$$\text{cenv}(\text{card } x) \text{ on } \{x \in \mathbb{R}^n \mid \|x\|_\infty \leq \kappa\} = \frac{1}{\kappa} \|x\|_1 \quad (1520)$$

$$\text{cenv}(\text{card } x) \text{ on } \{x \in \mathbb{R}_+^n \mid \|x\|_\infty \leq \kappa\} = \frac{1}{\kappa} \mathbf{1}^T x \quad (1521)$$

**7.2.2.2 Applying trace rank-heuristic to Problem 2**

Substituting rank envelope for rank function in Problem 2, for  $D \in \mathbb{EDM}^N$  ([confer \(1187\)](#))

$$\text{cenv rank}(-V_N^T D V_N) = \text{cenv rank}(-V D V) \propto -\text{tr}(V D V) \quad (1522)$$

and for desired affine dimension  $\rho \leq N - 1$  and nonnegative  $H$  [*sic*] we get a convex optimization problem

$$\begin{aligned} & \underset{D}{\text{minimize}} \quad \|\sqrt[D]{D} - H\|_{\text{F}}^2 \\ & \text{subject to} \quad -\text{tr}(V D V) \leq \kappa \rho \\ & \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1523)$$

where  $\kappa \in \mathbb{R}_+$  is a constant determined by cut-and-try. The equivalent semidefinite program makes  $\kappa$  variable: for nonnegative and symmetric  $H$

$$\begin{aligned} & \underset{D, Y, \kappa}{\text{minimize}} \quad \kappa \rho + 2 \text{tr}(V Y V) \\ & \text{subject to} \quad \begin{bmatrix} d_{ij} & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N - 1 \\ & \quad -\text{tr}(V D V) \leq \kappa \rho \\ & \quad Y \in \mathbb{S}_h^N \\ & \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1524)$$

<sup>7.19</sup> Provided  $f \neq +\infty$  and there exists an affine function  $h \leq f$  on  $\mathbb{R}^n$ , then the convex envelope is equal to the convex conjugate (the *Legendre-Fenchel transform*) of the convex conjugate of  $f$ ; *id est*, the conjugate-conjugate function  $f^{**}$ . [225, §E.1]

which is the same as (1513), the problem with no explicit constraint on affine dimension. As the present problem is stated, the desired affine dimension  $\rho$  yields to the variable scale factor  $\kappa$ ;  $\rho$  is effectively ignored.

Yet this result is an illuminant for problem (1513) and its equivalents (all the way back to (1506)): When the given measurement matrix  $H$  is nonnegative and symmetric, finding the closest EDM  $D$  (as in problem (1506), (1509), or (1513)) implicitly entails minimization of affine dimension (*confer* §5.8.4, §5.14.4). Those non-rank-constrained problems are each inherently equivalent to cenv(rank)-minimization problem (1524), in other words, and their optimal solutions are unique because of the strictly convex objective function in (1506).

### 7.2.2.3 Rank-heuristic insight

Minimization of affine dimension by use of this trace rank-heuristic (1522) tends to find a list configuration of least energy; rather, it tends to optimize compaction of the reconstruction by minimizing total distance. (1057) It is best used where some physical equilibrium implies such an energy minimization; *e.g.*, [393, §5].

For this Problem 2, the trace rank-heuristic arose naturally in the objective in terms of  $V$ . We observe:  $V$  (in contrast to  $V_N^T$ ) spreads energy over all available distances (§B.4.2 *no.20*, contrast *no.22*) although the rank function itself is insensitive to choice of auxiliary matrix.

Trace rank-heuristic (1518) is useless when a main diagonal is constrained to be constant. Such would be the case were optimization over an ellipope (§5.4.2.2.1), or when the diagonal represents a Boolean vector; *e.g.*, §4.2.3.1.1, §4.7.0.0.9.

### 7.2.2.4 Rank minimization heuristic beyond convex envelope

Fazel, Hindi, & Boyd [154] [449] [155] propose a rank heuristic more potent than trace (1519) for problems of rank minimization;

$$\text{rank } Y \leftarrow \log \det(Y + \varepsilon I) \quad (1525)$$

the concave surrogate function  $\log \det$  in place of quasiconcave  $\text{rank } Y$  (§2.9.2.9.2) when  $Y \in \mathbb{S}_+^n$  is variable and where  $\varepsilon$  is a small positive constant. They propose minimization of the surrogate by substituting a sequence comprising infima of a linearized surrogate about the current estimate  $Y_i$ ; *id est*, from the first-order Taylor series expansion about  $Y_i$  on some open interval of  $\|Y\|_2$  (§D.1.7)

$$\log \det(Y + \varepsilon I) \approx \log \det(Y_i + \varepsilon I) + \text{tr}((Y_i + \varepsilon I)^{-1}(Y - Y_i)) \quad (1526)$$

we make the surrogate sequence of infima over bounded convex feasible set  $\mathcal{C}$

$$\arg \inf_{Y \in \mathcal{C}} \text{rank } Y \leftarrow \lim_{i \rightarrow \infty} Y_{i+1} \quad (1527)$$

where, for  $i = 0 \dots$

$$Y_{i+1} = \arg \inf_{Y \in \mathcal{C}} \text{tr}((Y_i + \varepsilon I)^{-1}Y) \quad (1528)$$

a matrix analogue to the reweighting scheme disclosed in [234, §4.11.3]. Choosing  $Y_0 = I$ , the first step becomes equivalent to finding the infimum of  $\text{tr } Y$ ; the trace rank-heuristic (1519). The intuition underlying (1528) is the new term in the argument of trace; specifically,  $(Y_i + \varepsilon I)^{-1}$  weights  $Y$  so that relatively small eigenvalues of  $Y$  found by the infimum are made even smaller.

To see that, substitute the nonincreasingly ordered diagonalizations

$$\begin{aligned} Y_i + \varepsilon I &\triangleq Q(\Lambda + \varepsilon I)Q^T & \text{(a)} \\ Y &\triangleq U\Upsilon U^T & \text{(b)} \end{aligned} \quad (1529)$$

into (1528). Then from (1873) we have,

$$\begin{aligned} \inf_{\Upsilon \in U^* \cap U^*} \delta((\Lambda + \varepsilon I)^{-1})^T \delta(\Upsilon) &= \inf_{\Upsilon \in U^* \cap U^*} \inf_{R^T = R^{-1}} \text{tr}((\Lambda + \varepsilon I)^{-1} R^T \Upsilon R) \\ &\leq \inf_{Y \in \mathcal{C}} \text{tr}((Y_i + \varepsilon I)^{-1} Y) \end{aligned} \quad (1530)$$

where  $R \triangleq Q^T U$  in  $U$  on the set of orthogonal matrices is a bijection. The role of  $\varepsilon$  is, therefore, to limit maximum weight; the smallest entry on the main diagonal of  $\Upsilon$  gets the largest weight. ♦

#### 7.2.2.5 Applying log det rank-heuristic to Problem 2

When the log det rank-heuristic is inserted into Problem 2, problem (1524) becomes the problem sequence in  $i$

$$\begin{aligned} &\underset{D, Y, \kappa}{\text{minimize}} \quad \kappa \rho + 2 \text{tr}(VYV) \\ &\text{subject to} \quad \begin{bmatrix} d_{jl} & y_{jl} \\ y_{jl} & h_{jl}^2 \end{bmatrix} \succeq 0, \quad l > j = 1 \dots N-1 \\ &\quad -\text{tr}((-VD_i V + \varepsilon I)^{-1} V D V) \leq \kappa \rho \\ &\quad Y \in \mathbb{S}_h^N \\ &\quad D \in \mathbb{EDM}^N \end{aligned} \quad (1531)$$

where  $D_{i+1} \triangleq D^* \in \mathbb{EDM}^N$  and  $D_0 \triangleq \mathbf{1}\mathbf{1}^T - I$ .

#### 7.2.2.6 Tightening this log det rank-heuristic

Like the trace method, this log det technique for constraining rank offers no provision for meeting a predetermined upper bound  $\rho$ . Yet since eigenvalues are simply determined,  $\lambda(Y_i + \varepsilon I) = \delta(\Lambda + \varepsilon I)$ , we may certainly force selected weights to  $\varepsilon^{-1}$  by manipulating diagonalization (1529a). Empirically we find this sometimes leads to better results, although affine dimension of a solution cannot be guaranteed.

#### 7.2.2.7 Cumulative summary of rank heuristics

We have studied a perturbation method of rank reduction in §4.3 as well as the trace heuristic (convex envelope method §7.2.2.1.1) and log det heuristic in §7.2.2.4. There is another good contemporary method called LMIRank [317] based on alternating projection (§E.10).<sup>7.20</sup>

##### 7.2.2.7.1 Example. Unidimensional scaling.

We apply the convex iteration method from §4.5.1 to numerically solve an instance of Problem 2; a method empirically superior to the foregoing convex envelope and log det heuristics for rank regularization and enforcing affine dimension.

*Unidimensional scaling*, [117] a historically practical application of multidimensional scaling (§5.12), entails solution of an optimization problem having local minima whose

---

<sup>7.20</sup> that does not solve the *ball packing* problem presented in §5.4.2.2.6.

multiplicity varies as the factorial of point-list cardinality; geometrically, it means reconstructing a list constrained to lie in one affine dimension. Given nonnegative symmetric matrix  $H = [h_{ij}] \in \mathbb{S}^N \cap \mathbb{R}_+^{N \times N}$  (1511) whose entries  $h_{ij}$  are all known, the nonconvex problem in terms of point list is

$$\underset{\{x_i \in \mathbb{R}\}}{\text{minimize}} \sum_{i,j=1}^N (|x_i - x_j| - h_{ij})^2 \quad (1454)$$

called a *raw stress* problem [56, p.34] which has an implicit constraint on dimensional embedding of points  $\{x_i \in \mathbb{R}, i=1 \dots N\}$ . This problem has proven NP-hard; e.g, [80].

As always, we first transform variables to distance-square  $D \in \mathbb{S}_h^N$ ; so begin with convex problem (1513) on page 470

$$\begin{aligned} & \underset{D, Y}{\text{minimize}} \quad -\text{tr}(V(D - 2Y)V) \\ \text{subject to} \quad & \begin{bmatrix} d_{ij} & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N-1 \\ & Y \in \mathbb{S}_h^N \\ & D \in \mathbb{EDM}^N \\ & \text{rank } V_N^T DV_N = 1 \end{aligned} \quad (1532)$$

that becomes equivalent to (1454) by making explicit the constraint on affine dimension via rank. The iteration is formed by moving the dimensional constraint to the objective:

$$\begin{aligned} & \underset{D, Y}{\text{minimize}} \quad -\langle V(D - 2Y)V, I \rangle - w \langle V_N^T DV_N, W \rangle \\ \text{subject to} \quad & \begin{bmatrix} d_{ij} & y_{ij} \\ y_{ij} & h_{ij}^2 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N-1 \\ & Y \in \mathbb{S}_h^N \\ & D \in \mathbb{EDM}^N \end{aligned} \quad (1533)$$

where  $w$  ( $\approx 10$ ) is a positive scalar just large enough to make  $\langle V_N^T DV_N, W \rangle$  vanish to within some numerical precision, and where direction matrix  $W$  is an optimal solution to semidefinite program (1872a)

$$\begin{aligned} & \underset{W}{\text{minimize}} \quad -\langle V_N^T D^* V_N, W \rangle \\ \text{subject to} \quad & 0 \preceq W \preceq I \\ & \text{tr } W = N-1 \end{aligned} \quad (1534)$$

one of which is known in closed form. Semidefinite programs (1533) and (1534) are iterated until convergence in the sense defined on page 248. This iteration is not a projection method. (§4.5.1.1) Convex problem (1533) is neither a relaxation of unidimensional scaling problem (1532); instead, problem (1533) is a convex equivalent to (1532) at convergence of the iteration.

Jan de Leeuw provided us with some test data

$$H = \begin{bmatrix} 0.000000 & 5.235301 & 5.499274 & 6.404294 & 6.486829 & 6.263265 \\ 5.235301 & 0.000000 & 3.208028 & 5.840931 & 3.559010 & 5.353489 \\ 5.499274 & 3.208028 & 0.000000 & 5.679550 & 4.020339 & 5.239842 \\ 6.404294 & 5.840931 & 5.679550 & 0.000000 & 4.862884 & 4.543120 \\ 6.486829 & 3.559010 & 4.020339 & 4.862884 & 0.000000 & 4.618718 \\ 6.263265 & 5.353489 & 5.239842 & 4.543120 & 4.618718 & 0.000000 \end{bmatrix} \quad (1535)$$

and a globally optimal solution

$$\begin{aligned} X^* &= [-4.981494 \quad -2.121026 \quad -1.038738 \quad 4.555130 \quad 0.764096 \quad 2.822032] \\ &= [x_1^* \quad x_2^* \quad x_3^* \quad x_4^* \quad x_5^* \quad x_6^*] \end{aligned} \quad (1536)$$

found by searching 6! local minima of (1454) [117]. By iterating convex problems (1533) and (1534) about twenty times (initial  $W = \mathbf{0}$ ) we find global infimum 98.12812 to stress problem (1454), and by (1277) we find a corresponding one-dimensional point list that is a rigid transformation in  $\mathbb{R}$  of  $X^*$ .

Here we found the infimum to accuracy of the given data, but that ceases to hold as problem size increases. Because of machine numerical precision and an interior-point method of solution, we speculate, accuracy degrades quickly as problem size increases beyond this.  $\square$

### 7.3 Third prevalent problem: Projection on EDM cone in $d_{ij}$

*In summary, we find that the solution to problem [(1456.3) p.459] is difficult and depends on the dimension of the space as the geometry of the cone of EDMs becomes more complex.*

—Hayden, Wells, Liu, & Tarazaga, 1991 [211, §3]

Reformulating Problem 2 (p.468), in terms of EDM  $D$ , changes it considerably:

$$\left. \begin{array}{ll} \text{minimize}_D & \|D - H\|_F^2 \\ \text{subject to} & \text{rank } V_N^T D V_N \leq \rho \\ & D \in \mathbb{EDM}^N \end{array} \right\} \quad \text{Problem 3} \quad (1537)$$

This third prevalent proximity problem is a Euclidean projection of given matrix  $H$  on a generally nonconvex subset ( $\rho < N - 1$ ) of  $\partial \mathbb{EDM}^N$  the boundary of the convex cone of Euclidean distance matrices relative to subspace  $\mathbb{S}_h^N$  (Figure 163d). Because coordinates of projection are distance-square and  $H$  now presumably holds distance-square measurements, numerical solution to Problem 3 is generally different than that of Problem 2.

For the moment, we need make no assumptions regarding measurement matrix  $H$ .

#### 7.3.1 Convex case

$$\left. \begin{array}{ll} \text{minimize}_D & \|D - H\|_F^2 \\ \text{subject to} & D \in \mathbb{EDM}^N \end{array} \right\} \quad (1538)$$

When the rank constraint disappears (for  $\rho = N - 1$ ), this third problem becomes obviously convex because the feasible set is then the entire EDM cone and because the objective function

$$\|D - H\|_F^2 = \sum_{i,j} (d_{ij} - h_{ij})^2 \quad (1539)$$

is a strictly convex quadratic in  $D$ ;<sup>7.21</sup>

$$\begin{aligned} & \underset{D}{\text{minimize}} \quad \sum_{i,j} d_{ij}^2 - 2h_{ij}d_{ij} + h_{ij}^2 \\ & \text{subject to} \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1540)$$

Optimal solution  $D^*$  is therefore unique, as expected, for this simple projection on the EDM cone equivalent to (1455).

### 7.3.1.1 Equivalent semidefinite program, Problem 3, convex case

In the past, this convex problem was solved numerically by means of alternating projection. (Example 7.3.1.1.1) [176] [168] [211, §1] We translate (1540) to an equivalent semidefinite program because we have a good solver:

Assume the given measurement matrix  $H$  to be nonnegative and symmetric;<sup>7.22</sup>

$$H = [h_{ij}] \in \mathbb{S}^N \cap \mathbb{R}_+^{N \times N} \quad (1511)$$

We then propose: Problem (1540) is equivalent to the semidefinite program, for

$$\partial \triangleq [d_{ij}^2] = D \circ D \quad (1541)$$

a matrix of distance-square squared,

$$\begin{aligned} & \underset{\partial, D}{\text{minimize}} \quad -\text{tr}(V(\partial - 2H \circ D)V) \\ & \text{subject to} \quad \begin{bmatrix} \partial_{ij} & d_{ij} \\ d_{ij} & 1 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N-1 \\ & \quad D \in \mathbb{EDM}^N \\ & \quad \partial \in \mathbb{S}_h^N \end{aligned} \quad (1542)$$

where

$$\begin{bmatrix} \partial_{ij} & d_{ij} \\ d_{ij} & 1 \end{bmatrix} \succeq 0 \Leftrightarrow \partial_{ij} \geq d_{ij}^2 \quad (1543)$$

Symmetry of input  $H$  facilitates trace in the objective (§B.4.2 no.20), while its nonnegativity causes  $\partial_{ij} \rightarrow d_{ij}^2$  as optimization proceeds.

#### 7.3.1.1.1 Example. Alternating projection on nearest EDM.

By solving (1542) we confirm the result from an example given by Glunt, Hayden, Hong, & Wells [176, §6] who found analytical solution to convex optimization problem (1538) for particular cardinality  $N=3$  by using the alternating projection method of von Neumann (§E.10):

$$H = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 9 \\ 1 & 9 & 0 \end{bmatrix}, \quad D^* = \begin{bmatrix} 0 & \frac{19}{9} & \frac{19}{9} \\ \frac{19}{9} & 0 & \frac{76}{9} \\ \frac{19}{9} & \frac{76}{9} & 0 \end{bmatrix} \quad (1544)$$

---

<sup>7.21</sup>For nonzero  $Y \in \mathbb{S}_h^N$  and some open interval of  $t \in \mathbb{R}$  (§3.13.0.0.2, §D.2.3)

$$\frac{d^2}{dt^2} \| (D + tY) - H \|_F^2 = 2 \text{tr} Y^T Y > 0 \quad \blacklozenge$$

<sup>7.22</sup>If that  $H$  given has negative entries, then the technique of solution presented here becomes invalid. Projection of  $H$  on  $\mathcal{K}$  (1447) prior to application of this proposed technique, as explained in §7.0.1, is incorrect.

The problem (1538), of projecting  $H$  on the EDM cone, is transformed to an equivalent iterative sequence of projections on the two convex cones (1396) from §6.8.1.1. Utilizing projector (1399) in an ordinary alternating projection, input  $H$  goes to  $D^*$  with an accuracy of four decimal places in about 17 iterations. Affine dimension corresponding to this optimal solution is  $r = 1$ .

Obviation of semidefinite programming's computational expense is the principal advantage of this alternating projection technique.  $\square$

### 7.3.1.2 Schur-form semidefinite program, Problem 3 convex case

Semidefinite program (1542) can be reformulated by moving the objective function in

$$\begin{aligned} & \underset{D}{\text{minimize}} \quad \|D - H\|_F^2 \\ & \text{subject to} \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1538)$$

to the constraints. This makes an equivalent epigraph form of the problem: for any measurement matrix  $H$

$$\begin{aligned} & \underset{t \in \mathbb{R}, D}{\text{minimize}} \quad t \\ & \text{subject to} \quad \|D - H\|_F^2 \leq t \\ & \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1545)$$

We can transform this problem to an equivalent Schur-form semidefinite program; (§3.5.3)

$$\begin{aligned} & \underset{t \in \mathbb{R}, D}{\text{minimize}} \quad t \\ & \text{subject to} \quad \begin{bmatrix} tI & \text{vec}(D - H) \\ \text{vec}(D - H)^T & 1 \end{bmatrix} \succeq 0 \\ & \quad D \in \mathbb{EDM}^N \end{aligned} \quad (1546)$$

characterized by great sparsity and structure. The advantage of this SDP is lack of conditions on input  $H$ ; *e.g.*, negative entries would invalidate any solution provided by (1542). (§7.0.1.2)

### 7.3.1.3 Gram-form semidefinite program, Problem 3 convex case

Further, this problem statement may be equivalently written in terms of a Gram matrix via linear bijective (§5.6.1) EDM operator  $\mathbf{D}(G)$  (1045);

$$\begin{aligned} & \underset{G \in \mathbb{S}_c^N, t \in \mathbb{R}}{\text{minimize}} \quad t \\ & \text{subject to} \quad \begin{bmatrix} tI & \text{vec}(\mathbf{D}(G) - H) \\ \text{vec}(\mathbf{D}(G) - H)^T & 1 \end{bmatrix} \succeq 0 \\ & \quad G \succeq 0 \end{aligned} \quad (1547)$$

To include constraints on the list  $X \in \mathbb{R}^{n \times N}$ , we would rewrite this:

$$\begin{aligned} & \underset{G \in \mathbb{S}_c^N, t \in \mathbb{R}, X \in \mathbb{R}^{n \times N}}{\text{minimize}} \quad t \\ & \text{subject to} \quad \begin{bmatrix} tI & \text{vec}(\mathbf{D}(G) - H) \\ \text{vec}(\mathbf{D}(G) - H)^T & 1 \end{bmatrix} \succeq 0 \\ & \quad \begin{bmatrix} I & X \\ X^T & G \end{bmatrix} \succeq 0 \\ & \quad X \in \mathcal{C} \end{aligned} \quad (1548)$$

where  $\mathcal{C}$  is some abstract convex set. This technique is discussed in §5.4.2.2.8.

### 7.3.1.4 Dual interpretation, projection on EDM cone

From §E.9.1.2 we learn that projection on a convex cone has dual form. In the circumstance that  $\mathcal{K}$  is a convex cone and point  $x$  exists exterior to the cone or on its boundary, distance to the nearest point  $Px$  in  $\mathcal{K}$  is found as the optimal value of the objective

$$\begin{aligned} \|x - Px\| &= \underset{a}{\text{maximize}} \quad a^T x \\ \text{subject to} \quad \|a\| &\leq 1 \\ a &\in \mathcal{K}^\circ \end{aligned} \quad (2215)$$

where  $\mathcal{K}^\circ$  is the polar cone.

Applying this result to (1538), we get a convex optimization for any given symmetric matrix  $H$  exterior to or on the EDM cone boundary:

$$\begin{aligned} \underset{D}{\text{minimize}} \quad \|D - H\|_F^2 &\quad \underset{\substack{A^\circ \\ A^\circ \in \text{EDM}^{N^\circ}}}{\text{maximize}} \quad \langle A^\circ, H \rangle \\ \text{subject to} \quad D &\in \text{EDM}^N \quad \equiv \quad \text{subject to} \quad \|A^\circ\|_F \leq 1 \end{aligned} \quad (1549)$$

Then, from (2217), projection of  $H$  on cone  $\text{EDM}^N$  is

$$D^* = H - A^{\circ*} \langle A^{\circ*}, H \rangle \quad (1550)$$

Critchley proposed, instead, projection on the polar EDM cone in his 1980 thesis [97, p.113]: In that circumstance, by projection on the algebraic complement (§E.9.2.2.1),

$$D^* = A^* \langle A^*, H \rangle \quad (1551)$$

which is equal to (1550) when  $A^*$  solves

$$\begin{aligned} \underset{A}{\text{maximize}} \quad \langle A, H \rangle & \\ \text{subject to} \quad \|A\|_F &= 1 \\ A &\in \text{EDM}^N \end{aligned} \quad (1552)$$

This projection of symmetric  $H$  on polar cone  $\text{EDM}^{N^\circ}$  can be made a convex problem, of course, by relaxing the equality constraint ( $\|A\|_F \leq 1$ ).

### 7.3.2 Minimization of affine dimension in Problem 3

When desired affine dimension  $\rho$  is diminished, Problem 3 (1537) is difficult to solve [211, §3] because the feasible set in  $\mathbb{R}^{N(N-1)/2}$  loses convexity. By substituting rank envelope (1522) into Problem 3, then for any given  $H$  we get a convex problem

$$\begin{aligned} \underset{D}{\text{minimize}} \quad \|D - H\|_F^2 & \\ \text{subject to} \quad -\text{tr}(VDV) &\leq \kappa \rho \\ D &\in \text{EDM}^N \end{aligned} \quad (1553)$$

where  $\kappa \in \mathbb{R}_+$  is a constant determined by cut-and-try. Given  $\kappa$ , problem (1553) is a convex optimization having unique solution in any desired affine dimension  $\rho$ ; an approximation to Euclidean projection on that nonconvex subset of the EDM cone containing EDMs with corresponding affine dimension no greater than  $\rho$ .

The SDP equivalent to (1553) does not move  $\kappa$  into the variables as on page 472: for nonnegative symmetric input  $H$  and distance-square squared variable  $\partial$  as in (1541),

$$\begin{aligned} & \underset{\partial, D}{\text{minimize}} \quad -\text{tr}(V(\partial - 2H \circ D)V) \\ & \text{subject to} \quad \begin{bmatrix} \partial_{ij} & d_{ij} \\ d_{ij} & 1 \end{bmatrix} \succeq 0, \quad N \geq j > i = 1 \dots N-1 \\ & \quad -\text{tr}(VDV) \leq \kappa\rho \\ & \quad D \in \mathbb{EDM}^N \\ & \quad \partial \in \mathbb{S}_h^N \end{aligned} \tag{1554}$$

That means we will not see equivalence of this cenv(rank)-minimization problem to the non-rank-constrained problems (1540) and (1542) like we saw for its counterpart (1524) in Problem 2.

Another approach to affine dimension minimization is to project instead on the polar EDM cone; discussed in §6.8.1.5.

### 7.3.3 Constrained affine dimension, Problem 3

When one desires affine dimension diminished further below what can be achieved via cenv(rank)-minimization as in (1554), spectral projection can be considered a natural means in light of its successful application to projection on a rank  $\rho$  subset of a positive semidefinite cone in §7.1.4.

Yet it is wrong here to zero eigenvalues of  $-VDV$  or  $-VGV$  or a variant to reduce affine dimension, because that particular method comes from projection on a positive semidefinite cone (1483); zeroing those eigenvalues here in Problem 3 would place an elbow in the projection path (Figure 178) thereby producing a result that is necessarily suboptimal. Problem 3 is instead a projection on the EDM cone whose associated spectral cone is considerably different. (§5.11.2.3) Proper choice of spectral cone is demanded by diagonalization of that variable argument to the objective:

#### 7.3.3.1 Cayley-Menger form

We use Cayley-Menger composition of the Euclidean distance matrix to solve a problem that is the same as Problem 3 (1537): (§5.7.3.0.1)

$$\begin{aligned} & \underset{D}{\text{minimize}} \quad \left\| \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} - \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -H \end{bmatrix} \right\|_F^2 \\ & \text{subject to} \quad \text{rank} \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} \leq \rho + 2 \\ & \quad D \in \mathbb{EDM}^N \end{aligned} \tag{1555}$$

a projection of  $H$  on a generally nonconvex subset (when  $\rho < N-1$ ) of the Euclidean distance matrix cone boundary rel  $\partial \mathbb{EDM}^N$ ; *id est*, projection from the EDM cone interior or exterior on a subset of its relative boundary (§6.5, (1325)).

Rank of an optimal solution is intrinsically bounded above and below;

$$2 \leq \text{rank} \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D^* \end{bmatrix} \leq \rho + 2 \leq N + 1 \tag{1556}$$

Our proposed strategy for low-rank solution is projection on that subset of a spectral cone  $\lambda \left( \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -\mathbb{EDM}^N \end{bmatrix} \right)$  (§5.11.2.3) corresponding to affine dimension not in excess of that  $\rho$

desired; *id est*, spectral projection on

$$\begin{bmatrix} \mathbb{R}_+^{\rho+1} \\ \mathbf{0} \\ \mathbb{R}_- \end{bmatrix} \cap \partial\mathcal{H} \subset \mathbb{R}^{N+1} \quad (1557)$$

where

$$\partial\mathcal{H} = \{\lambda \in \mathbb{R}^{N+1} \mid \mathbf{1}^T \lambda = 0\} \quad (1257)$$

is a hyperplane through the origin. This pointed polyhedral cone (1557), to which membership subsumes the rank constraint, is not full-dimensional.

Given desired affine dimension  $0 \leq \rho \leq N-1$  and diagonalization (§A.5) of unknown EDM  $D$

$$\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix} \triangleq U\Upsilon U^T \in \mathbb{S}_h^{N+1} \quad (1558)$$

and given symmetric  $H$  in diagonalization

$$\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -H \end{bmatrix} \triangleq Q\Lambda Q^T \in \mathbb{S}^{N+1} \quad (1559)$$

having eigenvalues arranged in nonincreasing order, then by (1270) problem (1555) is equivalent to

$$\begin{aligned} & \underset{\Upsilon, R}{\text{minimize}} \quad \|\delta(\Upsilon) - \pi(\delta(R^T \Lambda R))\|^2 \\ & \text{subject to} \quad \delta(\Upsilon) \in \begin{bmatrix} \mathbb{R}_+^{\rho+1} \\ \mathbf{0} \\ \mathbb{R}_- \end{bmatrix} \cap \partial\mathcal{H} \\ & \quad \delta(QR\Upsilon R^T Q^T) = \mathbf{0} \\ & \quad R^{-1} = R^T \end{aligned} \quad (1560)$$

where  $\pi$  is the permutation operator from §7.1.3 arranging its vector argument in nonincreasing order,<sup>7.23</sup> where

$$R \triangleq Q^T U \in \mathbb{R}^{N+1 \times N+1} \quad (1561)$$

in  $U$  on the set of orthogonal matrices is a bijection, and where  $\partial\mathcal{H}$  insures one negative eigenvalue. HOLLOWNESS constraint  $\delta(QR\Upsilon R^T Q^T) = \mathbf{0}$  makes problem (1560) difficult by making the two variables dependent.

Our plan is to instead divide problem (1560) into two and then alternate their solution:

$$\begin{aligned} & \underset{\Upsilon}{\text{minimize}} \quad \|\delta(\Upsilon) - \pi(\delta(R^T \Lambda R))\|^2 \\ & \text{subject to} \quad \delta(\Upsilon) \in \begin{bmatrix} \mathbb{R}_+^{\rho+1} \\ \mathbf{0} \\ \mathbb{R}_- \end{bmatrix} \cap \partial\mathcal{H} \end{aligned} \quad (a) \quad (1562)$$

$$\begin{aligned} & \underset{R}{\text{minimize}} \quad \|R\Upsilon^* R^T - \Lambda\|_F^2 \\ & \text{subject to} \quad \delta(QR\Upsilon^* R^T Q^T) = \mathbf{0} \\ & \quad R^{-1} = R^T \end{aligned} \quad (b)$$

---

<sup>7.23</sup>Recall, any permutation matrix is an orthogonal matrix.

**Proof.** We justify disappearance of the hollowness constraint in convex optimization problem (1562a): From arguments in §7.1.3 with regard to permutation operator  $\pi$ , cone

membership constraint  $\delta(\Upsilon) \in \begin{bmatrix} \mathbb{R}_+^{\rho+1} \\ \mathbf{0} \\ \mathbb{R}_- \end{bmatrix} \cap \partial\mathcal{H}$  from (1562a) is equivalent to

$$\delta(\Upsilon) \in \begin{bmatrix} \mathbb{R}_+^{\rho+1} \\ \mathbf{0} \\ \mathbb{R}_- \end{bmatrix} \cap \partial\mathcal{H} \cap \mathcal{K}_{\mathcal{M}} \quad (1563)$$

where  $\mathcal{K}_{\mathcal{M}}$  is the monotone cone (§2.13.10.4.3). Membership of  $\delta(\Upsilon)$  to the *polyhedral cone of majorization* (Theorem A.1.2.0.1)

$$\mathcal{K}_{\lambda\delta}^* = \partial\mathcal{H} \cap \mathcal{K}_{\mathcal{M}+}^* \quad (1580)$$

where  $\mathcal{K}_{\mathcal{M}+}^*$  is the dual monotone nonnegative cone (§2.13.10.4.2), is a condition (in absence of a hollowness constraint) that would insure existence of a symmetric hollow matrix  $\begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & -D \end{bmatrix}$ . Curiously, intersection of this feasible superset  $\begin{bmatrix} \mathbb{R}_+^{\rho+1} \\ \mathbf{0} \\ \mathbb{R}_- \end{bmatrix} \cap \partial\mathcal{H} \cap \mathcal{K}_{\mathcal{M}}$  from (1563) with the cone of majorization  $\mathcal{K}_{\lambda\delta}^*$  is a benign operation; *id est*,

$$\partial\mathcal{H} \cap \mathcal{K}_{\mathcal{M}+}^* \cap \mathcal{K}_{\mathcal{M}} = \partial\mathcal{H} \cap \mathcal{K}_{\mathcal{M}} \quad (1564)$$

verifiable by observing conic dependencies (§2.10.3) among the aggregate of halfspace-description normals. The cone membership constraint in (1562a) therefore inherently insures existence of a symmetric hollow matrix. ♦

Optimization (1562b) would be a Procrustes problem (§C.4) were it not for the hollowness constraint. It is, instead, a minimization over the intersection of the nonconvex manifold of orthogonal matrices with another nonconvex set in variable  $R$  specified by the hollowness constraint. We solve problem (1562b) by a method introduced in §4.7.0.0.2: Define  $R = [r_1 \cdots r_{N+1}] \in \mathbb{R}^{N+1 \times N+1}$  and make the assignment

$$G = \begin{bmatrix} r_1 \\ \vdots \\ r_{N+1} \\ 1 \end{bmatrix} [r_1^T \cdots r_{N+1}^T \ 1] \in \mathbb{S}^{(N+1)^2+1} \\ = \begin{bmatrix} R_{11} & \cdots & R_{1,N+1} & r_1 \\ \vdots & \ddots & \vdots & \vdots \\ R_{1,N+1}^T & \cdots & R_{N+1,N+1} & r_{N+1} \end{bmatrix} \triangleq \begin{bmatrix} r_1 r_1^T & \cdots & r_1 r_{N+1}^T & r_1 \\ \vdots & \ddots & \vdots & \vdots \\ r_{N+1} r_1^T & \cdots & r_{N+1} r_{N+1}^T & r_{N+1} \\ r_1^T & \cdots & r_{N+1}^T & 1 \end{bmatrix} \quad (1565)$$

where  $R_{ij} \triangleq r_i r_j^T \in \mathbb{R}^{N+1 \times N+1}$  and  $\Upsilon_{ii}^* \in \mathbb{R}$ . Since  $R \Upsilon^* R^T = \sum_{i=1}^{N+1} \Upsilon_{ii}^* R_{ii}$ , then problem

(1562b) is equivalently expressed:

$$\begin{aligned}
 & \underset{R_{ii} \in \mathbb{S}, R_{ij}, r_i}{\text{minimize}} \quad \left\| \sum_{i=1}^{N+1} \Upsilon_{ii}^* R_{ii} - \Lambda \right\|_F^2 \\
 & \text{subject to} \quad \begin{aligned} \text{tr } R_{ii} &= 1, & i &= 1 \dots N+1 \\ \text{tr } R_{ij} &= 0, & i < j &= 2 \dots N+1 \end{aligned} \\
 & \quad G = \begin{bmatrix} R_{11} & \cdots & R_{1,N+1} & r_1 \\ \vdots & \ddots & & \vdots \\ R_{1,N+1}^T & \cdots & R_{N+1,N+1} & r_{N+1} \\ r_1^T & \cdots & r_{N+1}^T & 1 \end{bmatrix} (\succeq 0) \\
 & \quad \delta \left( Q \sum_{i=1}^{N+1} \Upsilon_{ii}^* R_{ii} Q^T \right) = \mathbf{0} \\
 & \quad \text{rank } G = 1
 \end{aligned} \tag{1566}$$

The rank constraint is regularized by method of convex iteration developed in §4.5. Problem (1566) is partitioned into two convex problems:

$$\begin{aligned}
 & \underset{R_{ij}, r_i}{\text{minimize}} \quad \left\| \sum_{i=1}^{N+1} \Upsilon_{ii}^* R_{ii} - \Lambda \right\|_F^2 + \langle G, W \rangle \\
 & \text{subject to} \quad \begin{aligned} \text{tr } R_{ii} &= 1, & i &= 1 \dots N+1 \\ \text{tr } R_{ij} &= 0, & i < j &= 2 \dots N+1 \end{aligned} \\
 & \quad G = \begin{bmatrix} R_{11} & \cdots & R_{1,N+1} & r_1 \\ \vdots & \ddots & & \vdots \\ R_{1,N+1}^T & \cdots & R_{N+1,N+1} & r_{N+1} \\ r_1^T & \cdots & r_{N+1}^T & 1 \end{bmatrix} \succeq 0 \\
 & \quad \delta \left( Q \sum_{i=1}^{N+1} \Upsilon_{ii}^* R_{ii} Q^T \right) = \mathbf{0}
 \end{aligned} \tag{1567}$$

and

$$\begin{aligned}
 & \underset{W \in \mathbb{S}^{(N+1)^2+1}}{\text{minimize}} \quad \langle G^*, W \rangle \\
 & \text{subject to} \quad \begin{aligned} 0 &\preceq W \preceq I \\ \text{tr } W &= (N+1)^2 \end{aligned}
 \end{aligned} \tag{1568}$$

then alternated with convex problem (1562a) until a rank-1  $G$  matrix is found and the objective of (1562a) is minimized.<sup>7.24</sup> An optimal solution to (1568) is known in closed form (p.533).

## 7.4 Conclusion

The importance and application of solving rank- or cardinality-constrained problems are enormous, a conclusion generally accepted *gratis* by the mathematics and engineering communities. Rank-constrained semidefinite programs arise in many vital feedback and control problems [196], optics [82] (Figure 136), and communications [161] [284] (Figure 114). *For example, one might be interested in the minimal order dynamic output feedback which stabilizes a given linear time invariant plant (this problem is considered to be among the most important open problems in control).* – [294] Rank and cardinality constraints also arise naturally in combinatorial optimization (§4.7.0.0.11, Figure 125), and find application to facial recognition (Figure 6), cartography (Figure 159), latent

<sup>7.24</sup>The hollowness constraint in (1567) may cause numerical instability; in that case, it may be moved to the objective within an added weighted norm. Conditions for termination of the iteration would then comprise a vanishing norm of hollowness.

semantic indexing [259], sparse or low-rank matrix completion for preference models and collaborative filtering, multidimensional scaling or principal component analysis (§5.12), medical imaging (Figure 120), digital filter design with time domain constraints [420], molecular conformation (Figure 152), sensor-network localization and wireless location (Figure 97), *etcetera*.

There has been little progress in spectral projection since the discovery by Eckart & Young in 1936 [147] leading to a formula for projection on a rank  $\rho$  subset of a positive semidefinite cone (§2.9.2.1). [166] The only closed-form spectral method presently available for solving proximity problems, having a constraint on rank, is based on their discovery (Problem 1, §7.1, §5.13).

- One popular recourse is intentional misapplication of Eckart & Young's result by introducing spectral projection on a positive semidefinite cone into Problem 3 via  $\mathbf{D}(G)$  (1045), for example. [80] Since Problem 3 instead demands spectral projection on the EDM cone, any solution acquired that way is necessarily suboptimal.
- A second recourse is problem redesign: A presupposition to all proximity problems in this chapter is that matrix  $H$  is given. We considered  $H$  having various properties such as nonnegativity, symmetry, hollowness, or lack thereof. It was assumed that if  $H$  did not already belong to the EDM cone, then we wanted an EDM closest to  $H$  in some sense; *id est*, input-matrix  $H$  was assumed corrupted somehow. For practical problems, it withstands reason that such a proximity problem could instead be reformulated so that some or all entries of  $H$  were unknown but bounded above and below by known limits; the norm objective is thereby eliminated as in the development beginning on page 259. That particular redesign (*the art*, p.8), in terms of the Gram-matrix bridge between point list  $X$  and EDM  $D$ , at once encompasses proximity and completion problems.
- A third recourse is to apply the method of convex iteration as we did in §7.2.2.7.1. This technique is applicable to any semidefinite problem requiring a rank constraint; it places a regularization term in the objective that enforces the rank constraint.

# Appendix A

## Linear algebra

### A.1 Main-diagonal $\delta$ operator, $\lambda$ , $\text{tr}$ , $\text{vec}$ , $\circ$ , $\otimes$

We introduce notation  $\delta$  denoting the main-diagonal linear selfadjoint operator. When linear function  $\delta$  operates on a square matrix  $A \in \mathbb{R}^{N \times N}$ ,  $\delta(A)$  returns a vector composed of all the entries from the main diagonal in the natural order;

$$\delta(A) \in \mathbb{R}^N \quad (1569)$$

Operating on a vector  $y \in \mathbb{R}^N$ ,  $\delta$  naturally returns a diagonal matrix;

$$\delta(y) \in \mathbb{S}^N \quad (1570)$$

Operating recursively on a vector  $\Lambda \in \mathbb{R}^N$  or diagonal matrix  $\Lambda \in \mathbb{S}^N$ ,  $\delta(\delta(\Lambda))$  returns  $\Lambda$  itself;

$$\delta^2(\Lambda) \equiv \delta(\delta(\Lambda)) \triangleq \Lambda \quad (1571)$$

Defined this way [254, §3.10, §9.5-1], [A.1](#) main-diagonal linear operator  $\delta$  is *selfadjoint*; *videlicet*, ([§2.2](#))

$$\delta(A)^T y = \langle \delta(A), y \rangle = \langle A, \delta(y) \rangle = \text{tr}(A^T \delta(y)) \quad (1572)$$

#### A.1.1 Identities

This  $\delta$  notation is efficient and unambiguous as illustrated in the following examples where:  $A \circ B$  denotes Hadamard's (commutative) product of matrices of like size [228] [181, §1.1.4] ([§D.1.2.2](#)),  $A \otimes B$  denotes Kronecker product [190] ([§D.1.2.1](#)),  $y$  is a vector,  $X$  a matrix,  $e_i$  the  $i^{\text{th}}$  member of the standard basis for  $\mathbb{R}^n$ ,  $\mathbb{S}_h^N$  the symmetric hollow subspace,  $\sigma(A)$  a vector of (nonincreasingly) ordered singular values of matrix  $A$ , and  $\lambda(A)$  denotes a vector of nonincreasingly ordered eigenvalues:

1.  $\delta(A) = \delta(A^T)$
2.  $\text{tr}(A) = \text{tr}(A^T) = \delta(A)^T \mathbf{1} = \langle I, A \rangle$
3.  $\delta(cA) = c \delta(A)$   $c \in \mathbb{R}$
4.  $\text{tr}(cA) = c \text{tr}(A) = c \mathbf{1}^T \lambda(A)$   $c \in \mathbb{R}$

---

[A.1](#) Linear operator  $T : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{M \times N}$  is selfadjoint when,  $\forall X_1, X_2 \in \mathbb{R}^{m \times n}$

$$\langle T(X_1), X_2 \rangle = \langle X_1, T(X_2) \rangle$$

5.  $\text{vec}(cA) = c \text{ vec}(A)$   $c \in \mathbb{R}$

6.  $A \circ B = B \circ A, \quad A \circ cB = cA \circ B$   $c \in \mathbb{R}$

7.  $A \otimes B \neq B \otimes A, \quad A \otimes cB = cA \otimes B$   $c \in \mathbb{R}$

8. There exist permutation matrices  $\Xi_1$  and  $\Xi_2$  [190, p.28]

$$A \otimes B = \Xi_1(B \otimes A) \Xi_2 \quad (1573)$$

9.  $\delta(A + B) = \delta(A) + \delta(B)$

10.  $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$

11.  $\text{vec}(A + B) = \text{vec}(A) + \text{vec}(B)$

12.  $(A + B) \circ C = A \circ C + B \circ C$   
 $A \circ (B + C) = A \circ B + A \circ C$

13.  $(A + B) \otimes C = A \otimes C + B \otimes C$   
 $A \otimes (B + C) = A \otimes B + A \otimes C$

14.  $\text{sgn}(c) \lambda(|c|A) = c \lambda(A)$   $c \in \mathbb{R}$

15.  $\text{sgn}(c) \sigma(|c|A) = c \sigma(A)$   $c \in \mathbb{R}$

16.  $\text{tr}(c\sqrt{A^T A}) = c \text{tr}\sqrt{A^T A} = c \mathbf{1}^T \sigma(A)$   $c \in \mathbb{R}$

17.  $\pi(\delta(A)) = \lambda(I \circ A)$  where  $\pi$  is presorting function

18.  $\delta(AB) = (A \circ B^T) \mathbf{1} = (B^T \circ A) \mathbf{1}, \quad \delta(AB)^T = \mathbf{1}^T (A^T \circ B) = \mathbf{1}^T (B \circ A^T)$

19.  $\delta(uv^T) = \begin{bmatrix} u_1 v_1 \\ \vdots \\ u_N v_N \end{bmatrix} = u \circ v,$   $u, v \in \mathbb{R}^N$

20.  $\langle \delta(uv^T), w \rangle = \langle u, \delta(wv^T) \rangle = \langle u \circ v, w \rangle = \langle u, w \circ v \rangle,$   $u, v, w \in \mathbb{R}^N$

21.  $\text{tr}(A^T B) = \text{tr}(AB^T) = \text{tr}(BA^T) = \text{tr}(B^T A) = \text{vec}(A)^T \text{vec } B$   
 $= \mathbf{1}^T (A \circ B) \mathbf{1} = \mathbf{1}^T \delta(AB^T) = \delta(A^T B)^T \mathbf{1} = \delta(BA^T)^T \mathbf{1} = \delta(B^T A)^T \mathbf{1}$

22.  $D = [d_{ij}] \in \mathbb{S}_h^N, \quad H = [h_{ij}] \in \mathbb{S}_h^N, \quad V = I - \frac{1}{N} \mathbf{1} \mathbf{1}^T \in \mathbb{S}^N$  (confer §B.4.2 no.20)

$$N \text{tr}(-V(D \circ H)V) = \text{tr}(D^T H) = \mathbf{1}^T (D \circ H) \mathbf{1} = \text{tr}(\mathbf{1} \mathbf{1}^T (D \circ H)) = \sum_{i,j} d_{ij} h_{ij}$$

23.  $\text{tr}(\Lambda A) = \delta(\Lambda)^T \delta(A), \quad \delta^2(\Lambda) \triangleq \Lambda \in \mathbb{S}^N$

24.  $y^T B \delta(A) = \text{tr}(B \delta(A) y^T) = \text{tr}(\delta(B^T y) A) = \text{tr}(A \delta(B^T y))$   
 $= \delta(A)^T B^T y = \text{tr}(y \delta(A)^T B^T) = \text{tr}(A^T \delta(B^T y)) = \text{tr}(\delta(B^T y) A^T)$

25.  $\delta^2(A^T A) = \sum_i e_i e_i^T A^T A e_i e_i^T$

26.  $\delta(\delta(A) \mathbf{1}^T) = \delta(\mathbf{1} \delta(A)^T) = \delta(A)$

27.  $\delta(A \mathbf{1}) \mathbf{1} = \delta(A \mathbf{1} \mathbf{1}^T) = A \mathbf{1}, \quad \delta(y) \mathbf{1} = \delta(y \mathbf{1}^T) = y$

28.  $\delta(I \mathbf{1}) = \delta(\mathbf{1}) = I$

29.  $\delta(e_i e_j^T \mathbf{1}) = \delta(e_i) = e_i e_i^T$

30. For  $\zeta = [\zeta_i] \in \mathbb{R}^k$  and  $x = [x_i] \in \mathbb{R}^k$ ,  $\sum_i \zeta_i / x_i = \zeta^T \delta(x)^{-1} \mathbf{1}$
31.  $\text{vec}(A \circ B) = \text{vec}(A) \circ \text{vec}(B) = \delta(\text{vec } A) \text{vec}(B)$  (42) (1963)  
 $= \text{vec}(B) \circ \text{vec}(A) = \delta(\text{vec } B) \text{vec}(A)$
32.  $\text{tr}(A^T X A) = \sum_{i=1}^n A(:, i)^T X A(:, i)$ ,  $A \in \mathbb{R}^{m \times n}$
33.  $\text{vec}(A X B) = (B^T \otimes A) \text{vec } X$  (not  $H$ )
34.  $\text{vec}(B X A) = (A^T \otimes B) \text{vec } X$
35.  $A X + X B = C \Leftrightarrow (I \otimes A + B^T \otimes I) \text{vec } X = \text{vec } C$  (Lyapunov) [333, §5.1.10]
36.  $\text{tr}(A X B X^T) = \text{vec}(X)^T \text{vec}(A X B) = \text{vec}(X)^T (B^T \otimes A) \text{vec } X$  [190]  
 $= \delta(\text{vec}(X) \text{vec}(X)^T (B^T \otimes A))^T \mathbf{1}$
37.  $\text{tr}(A X^T B X) = \text{vec}(X)^T \text{vec}(B X A) = \text{vec}(X)^T (A^T \otimes B) \text{vec } X$   
 $= \delta(\text{vec}(X) \text{vec}(X)^T (A^T \otimes B))^T \mathbf{1}$
38.  $a^T X B X^T c = \text{vec}(X)^T (B \otimes a c^T) \text{vec } X = \text{vec}(X)^T (B^T \otimes c a^T) \text{vec } X$  [333, §10.2.2]
39.  $a^T X^T B X c = \text{vec}(X)^T (a c^T \otimes B) \text{vec } X = \text{vec}(X)^T (c a^T \otimes B^T) \text{vec } X$

40. For any permutation matrix  $\Xi$  and dimensionally compatible vector  $y$  or matrix  $A$

$$\delta(\Xi y) = \Xi \delta(y) \Xi^T \quad (1574)$$

$$\delta(\Xi A \Xi^T) = \Xi \delta(A) \Xi^T \quad (1575)$$

So given any permutation matrix  $\Xi$  and any dimensionally compatible matrix  $B$ , for example,

$$\delta^2(B) = \Xi \delta^2(\Xi^T B \Xi) \Xi^T \quad (1576)$$

41.  $A \otimes 1 = 1 \otimes A = A$
42.  $A \otimes (B \otimes C) = (A \otimes B) \otimes C$
43.  $(A \otimes B)(D \otimes E) = AD \otimes BE$   
 $(A \otimes B \otimes C)(D \otimes E \otimes F) = AD \otimes BE \otimes CF$
44. For  $a$  a vector,  $(a \otimes B) = (a \otimes I)B$
45. For  $b^T$  a row vector,  $(A \otimes b^T) = A(I \otimes b^T)$
46.  $(A \otimes B)^T = A^T \otimes B^T$ ,  $(A \circ B)^T = A^T \circ B^T$
47.  $(AB)^{-1} = B^{-1}A^{-1}$ ,  $(AB)^T = B^T A^T$ ,  $(AB)^{-T} = A^{-T}B^{-T}$  (confer p.621)
48.  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$
49.  $(A \otimes B)^\dagger = A^\dagger \otimes B^\dagger$  (2054)
50.  $\delta(A \otimes B) = \delta(A) \otimes \delta(B)$
51.  $\text{tr}(A \otimes B) = \text{tr } A \text{tr } B = \text{tr}(B \otimes A)$  [262, cor.13.13]
52. For  $A \in \mathbb{R}^{m \times m}$ ,  $B \in \mathbb{R}^{n \times n}$ ,  $\det(A \otimes B) = \det^n(A) \det^m(B) = \det(B \otimes A)$
53. For eigenvalues  $\lambda(A) \in \mathbb{C}^n$  and eigenvectors  $v(A) \in \mathbb{C}^{n \times n}$   $\ni A = v\delta(\lambda)v^{-1} \in \mathbb{R}^{n \times n}$
- $\lambda(A \otimes B) = \lambda(A) \otimes \lambda(B)$ ,  $v(A \otimes B) = v(A) \otimes v(B)$  (1577)
54.  $\text{rank}(A \otimes B) = \text{rank}(A) \text{rank}(B) = \text{rank}(B \otimes A)$  [262, cor.13.11]

### A.1.2 Majorization

**A.1.2.0.1 Theorem.** (Schur) *Majorization.* [455, §7.4] [228, §4.3] [229, §5.5]  
 Let  $\lambda \in \mathbb{R}^N$  denote a given vector of eigenvalues and let  $\delta \in \mathbb{R}^N$  denote a given vector of main diagonal entries, both arranged in nonincreasing order. Then

$$\exists A \in \mathbb{S}^N \ni \lambda(A) = \lambda \text{ and } \delta(A) = \delta \Leftarrow \lambda - \delta \in \mathcal{K}_{\lambda\delta}^* \quad (1578)$$

and conversely

$$A \in \mathbb{S}^N \Rightarrow \lambda(A) - \delta(A) \in \mathcal{K}_{\lambda\delta}^* \quad (1579)$$

The difference belongs to the pointed polyhedral cone of majorization (not a full-dimensional cone, *confer*(318))

$$\mathcal{K}_{\lambda\delta}^* \triangleq \mathcal{K}_{\mathcal{M}^+}^* \cap \{\zeta \mathbf{1} \mid \zeta \in \mathbb{R}\}^* \quad (1580)$$

where  $\mathcal{K}_{\mathcal{M}^+}^*$  is the dual monotone nonnegative cone (440), and where the dual of the line is a hyperplane;  $\partial \mathcal{H} = \{\zeta \mathbf{1} \mid \zeta \in \mathbb{R}\}^* = \mathbf{1}^\perp$ .  $\diamond$

Majorization cone  $\mathcal{K}_{\lambda\delta}^*$  is naturally consequent to the definition of majorization; *id est*, vector  $y \in \mathbb{R}^N$  majorizes vector  $x$  if and only if

$$\sum_{i=1}^k x_i \leq \sum_{i=1}^k y_i \quad \forall 1 \leq k \leq N \quad (1581)$$

and

$$\mathbf{1}^T x = \mathbf{1}^T y \quad (1582)$$

Under these circumstances, rather, vector  $x$  is majorized by vector  $y$ .

In the particular circumstance  $\delta(A) = \mathbf{0}$  we get:

**A.1.2.0.2 Corollary.** *Symmetric hollow majorization.*

Let  $\lambda \in \mathbb{R}^N$  denote a given vector of eigenvalues arranged in nonincreasing order. Then

$$\exists A \in \mathbb{S}_h^N \ni \lambda(A) = \lambda \Leftarrow \lambda \in \mathcal{K}_{\lambda\delta}^* \quad (1583)$$

and conversely

$$A \in \mathbb{S}_h^N \Rightarrow \lambda(A) \in \mathcal{K}_{\lambda\delta}^* \quad (1584)$$

where  $\mathcal{K}_{\lambda\delta}^*$  is defined in (1580).  $\diamond$

## A.2 Semidefiniteness: domain of test

The most fundamental necessary, sufficient, and definitive test for positive semidefiniteness of matrix  $A \in \mathbb{R}^{n \times n}$  is: [229, §1]

$$x^T A x \geq 0 \text{ for each and every } x \in \mathbb{R}^n \text{ such that } \|x\| = 1 \quad (1585)$$

Traditionally, authors demand evaluation over broader domain; namely, over all  $x \in \mathbb{R}^n$  which is sufficient but unnecessary. Indeed, that standard textbook requirement is far over-reaching because if  $x^T A x$  is nonnegative for particular  $x = x_p$ , then it is nonnegative for any  $\alpha x_p$  where  $\alpha \in \mathbb{R}$ . Thus, only normalized  $x$  in  $\mathbb{R}^n$  need be evaluated.

Many authors add the further requirement that the domain be complex; the broadest domain. By so doing, only *Hermitian matrices* ( $A^H = A$  where superscript  $H$  denotes conjugate transpose)<sup>A.2</sup> are admitted to the set of positive semidefinite matrices (1588); an unnecessary prohibitive condition.

---

<sup>A.2</sup>Hermitian symmetry is the complex analogue; the real part of a Hermitian matrix is symmetric while its imaginary part is antisymmetric. A Hermitian matrix has real eigenvalues and real main diagonal.

### A.2.1 Symmetry *versus* semidefiniteness

We call (1585) *the most fundamental test* of positive semidefiniteness. Yet some authors instead say, for real  $A$  and complex domain  $\{x \in \mathbb{C}^n\}$ , the complex test  $x^H A x \geq 0$  is most fundamental. That complex broadening of the domain of test causes nonsymmetric real matrices to be excluded from the set of positive semidefinite matrices. Yet admitting nonsymmetric real matrices or not is a matter of preference<sup>A.3</sup> unless that complex test is adopted, as we shall now explain.

Any real square matrix  $A$  has a representation in terms of its symmetric and antisymmetric parts; *id est*,

$$A = \frac{(A + A^T)}{2} + \frac{(A - A^T)}{2} \quad (54)$$

Because, for all real  $A$ , the antisymmetric part vanishes under real test,

$$x^T \frac{(A - A^T)}{2} x = 0 \quad (1586)$$

only the symmetric part of  $A$ ,  $(A + A^T)/2$ , has a role determining positive semidefiniteness. Hence the oft-made presumption that only symmetric matrices may be positive semidefinite is, of course, erroneous under (1585). Because eigenvalue-signs of a symmetric matrix translate unequivocally to its semidefiniteness, the eigenvalues that determine semidefiniteness are always those of the *symmetrized* matrix. (§A.3) For that reason, and because symmetric (or Hermitian) matrices must have real eigenvalues, the convention adopted in the literature is that semidefinite matrices are synonymous with symmetric semidefinite matrices. Certainly misleading under (1585), that presumption is typically bolstered with compelling examples from the physical sciences where symmetric matrices occur within the mathematical exposition of natural phenomena.<sup>A.4</sup>

Perhaps a better explanation of this pervasive presumption of symmetry comes from Horn & Johnson [228, §7.1] whose perspective<sup>A.5</sup> is the complex matrix, thus necessitating the complex domain of test throughout their treatise. They explain, if  $A \in \mathbb{C}^{n \times n}$

*... and if  $x^H A x$  is real for all  $x \in \mathbb{C}^n$ , then  $A$  is Hermitian. Thus, the assumption that  $A$  is Hermitian is not necessary in the definition of positive definiteness. It is customary, however.*

Their comment is best explained by noting, the real part of  $x^H A x$  comes from the Hermitian part  $(A + A^H)/2$  of  $A$ ;

$$\operatorname{re}(x^H A x) = x^H \frac{A + A^H}{2} x \quad (1587)$$

rather,

$$x^H A x \in \mathbb{R} \Leftrightarrow A^H = A \quad (1588)$$

because the imaginary part of  $x^H A x$  comes from the anti-Hermitian part  $(A - A^H)/2$ ;

$$\operatorname{im}(x^H A x) = x^H \frac{A - A^H}{2} x \quad (1589)$$

that vanishes for nonzero  $x$  if and only if  $A = A^H$ . So the Hermitian symmetry assumption is unnecessary, according to Horn & Johnson, not because nonHermitian matrices could

<sup>A.3</sup>Golub & Van Loan [181, §4.2.2], for example, admit nonsymmetric real matrices.

<sup>A.4</sup>e.g, [159, §52] [348, §2.1]. Symmetric matrices are certainly pervasive in our chosen subject as well.

<sup>A.5</sup>A totally complex perspective is not necessarily more advantageous. The positive semidefinite cone, for example, is not selfdual (§2.13.6) in the ambient space of Hermitian matrices. [221, §II]

be regarded positive semidefinite, rather because nonHermitian (includes nonsymmetric real) matrices are not comparable on the real line under  $x^H Ax$ . Yet that complex edifice is dismantled in the test of real matrices (1585) because the domain of test is no longer necessarily complex; meaning,  $x^T Ax$  will certainly always be real, regardless of symmetry, and so real  $A$  will always be comparable.

In summary, if we limit the domain of test to all  $x$  in  $\mathbb{R}^n$  as in (1585), then nonsymmetric real matrices are admitted to the realm of semidefinite matrices because they become comparable on the real line. One important exception occurs for rank-one matrices  $\Psi = uv^T$  where  $u$  and  $v$  are real vectors:  $\Psi$  is positive semidefinite if and only if  $\Psi = uu^T$ . (§A.3.1.0.7)

We might choose to expand the domain of test to all  $x$  in  $\mathbb{C}^n$  so that only symmetric matrices would be comparable. An alternative to expanding domain of test is to assume all matrices of interest to be symmetric; that is commonly done, hence the synonymous relationship with semidefinite matrices.

#### A.2.1.0.1 Example. Nonsymmetric positive definite product.

Horn & Johnson assert and Zhang [455, §6.2, §3.2] agree:

*If  $A, B \in \mathbb{C}^{n \times n}$  are positive definite, then we know that the product  $AB$  is positive definite if and only if  $AB$  is Hermitian.* —[228, §7.6 prob.10]

Implicitly in their statement,  $A$  and  $B$  are assumed individually Hermitian and the domain of test is assumed complex. We prove the assertion to be false for real matrices under (1585) that adopts a real domain of test.

Proof is by counterexample:

$$A^T = A = \begin{bmatrix} 3 & 0 & -1 & 0 \\ 0 & 5 & 1 & 0 \\ -1 & 1 & 4 & 1 \\ 0 & 0 & 1 & 4 \end{bmatrix}, \quad \lambda(A) = \begin{bmatrix} 5.9 \\ 4.5 \\ 3.4 \\ 2.0 \end{bmatrix} \quad (1590)$$

$$B^T = B = \begin{bmatrix} 4 & 4 & -1 & -1 \\ 4 & 5 & 0 & 0 \\ -1 & 0 & 5 & 1 \\ -1 & 0 & 1 & 4 \end{bmatrix}, \quad \lambda(B) = \begin{bmatrix} 8.8 \\ 5.5 \\ 3.3 \\ 0.24 \end{bmatrix} \quad (1591)$$

$$(AB)^T \neq AB = \begin{bmatrix} 13 & 12 & -8 & -4 \\ 19 & 25 & 5 & 1 \\ -5 & 1 & 22 & 9 \\ -5 & 0 & 9 & 17 \end{bmatrix}, \quad \lambda(AB) = \begin{bmatrix} 36. \\ 29. \\ 10. \\ 0.72 \end{bmatrix} \quad (1592)$$

$$\frac{1}{2}(AB + (AB)^T) = \begin{bmatrix} 13 & 15.5 & -6.5 & -4.5 \\ 15.5 & 25 & 3 & 0.5 \\ -6.5 & 3 & 22 & 9 \\ -4.5 & 0.5 & 9 & 17 \end{bmatrix}, \quad \lambda\left(\frac{1}{2}(AB + (AB)^T)\right) = \begin{bmatrix} 36. \\ 30. \\ 10. \\ 0.014 \end{bmatrix} \quad (1593)$$

Whenever  $A \in \mathbb{S}_+^n$  and  $B \in \mathbb{S}_+^n$ , then  $\lambda(AB) = \lambda(\sqrt{A}B\sqrt{A})$  will always be a nonnegative vector by (1617) and Corollary A.3.1.0.5. Yet positive definiteness of product  $AB$  is certified instead by the nonnegative eigenvalues  $\lambda\left(\frac{1}{2}(AB + (AB)^T)\right)$  in (1593) (§A.3.1.0.1) despite the fact that  $AB$  is not symmetric. <sup>A.6</sup> ♦

Horn & Johnson and Zhang resolve this anomaly by choosing to exclude nonsymmetric matrices and products; they do so by expanding the domain of test to  $\mathbb{C}^n$ . □

<sup>A.6</sup>It is a little more difficult to find a counterexample in  $\mathbb{R}^{2 \times 2}$  or  $\mathbb{R}^{3 \times 3}$ ; which may have served to advance any confusion.

### A.3 Proper statements of positive semidefiniteness

Unlike Horn & Johnson and Zhang, we never adopt a complex domain of test with real matrices. So motivated is our consideration of proper statements of positive semidefiniteness under real domain of test. This restriction, ironically, complicates the facts when compared to corresponding statements for the complex case (found elsewhere [228] [455]).

We state several fundamental facts regarding positive semidefiniteness of real matrix  $A$  and the product  $AB$  and sum  $A+B$  of real matrices under fundamental real test (1585); a few require proof as they depart from the standard texts, while those remaining are well established or obvious.

#### A.3.0.0.1 Theorem. Positive (semi)definite matrix.

$A \in \mathbb{S}^M$  is positive semidefinite if and only if for each and every vector  $x \in \mathbb{R}^M$  of unit norm,  $\|x\| = 1$ ,<sup>A.7</sup> we have  $x^T A x \geq 0$  (1585);

$$A \succeq 0 \Leftrightarrow \text{tr}(xx^T A) = x^T A x \geq 0 \quad \forall xx^T \quad (1594)$$

Matrix  $A \in \mathbb{S}^M$  is positive definite if and only if for each and every  $\|x\| = 1$  we have  $x^T A x > 0$ ;

$$A \succ 0 \Leftrightarrow \text{tr}(xx^T A) = x^T A x > 0 \quad \forall xx^T, \quad xx^T \neq \mathbf{0} \quad (1595)$$

◊

**Proof.** Statements (1594) and (1595) are each a particular instance of dual generalized inequalities (§2.13.2) with respect to the positive semidefinite cone; *videlicet*, [401]

$$\begin{aligned} A \succeq 0 &\Leftrightarrow \langle xx^T, A \rangle \geq 0 \quad \forall xx^T (\succeq 0) \\ A \succ 0 &\Leftrightarrow \langle xx^T, A \rangle > 0 \quad \forall xx^T (\succeq 0), \quad xx^T \neq \mathbf{0} \end{aligned} \quad (1596)$$

This says: positive semidefinite matrix  $A$  must belong to the normal side of every hyperplane whose normal is an extreme direction of the positive semidefinite cone. Relations (1594) and (1595) remain true when  $xx^T$  is replaced with “for each and every” positive semidefinite matrix  $X \in \mathbb{S}_+^M$  (§2.13.6) of unit norm,  $\|X\| = 1$ , as in

$$\begin{aligned} A \succeq 0 &\Leftrightarrow \text{tr}(XA) \geq 0 \quad \forall X \in \mathbb{S}_+^M \\ A \succ 0 &\Leftrightarrow \text{tr}(XA) > 0 \quad \forall X \in \mathbb{S}_+^M, \quad X \neq \mathbf{0} \end{aligned} \quad (1597)$$

But that condition is more than what is necessary. By the *discretized membership theorem* in §2.13.4.2.1, the extreme directions  $xx^T$  of the positive semidefinite cone constitute a minimal set of generators necessary and sufficient for discretization of dual generalized inequalities (1597) certifying membership to that cone. ♦

#### A.3.1 Semidefiniteness, eigenvalues, nonsymmetric

When  $A \in \mathbb{R}^{n \times n}$ , let  $\lambda(\frac{1}{2}(A + A^T)) \in \mathbb{R}^n$  denote eigenvalues of the symmetrized matrix<sup>A.8</sup> arranged in nonincreasing order.

---

<sup>A.7</sup>The traditional condition requiring all  $x \in \mathbb{R}^M$  for defining positive (semi)definiteness is actually more than what is necessary. The set of norm-1 vectors is necessary and sufficient to establish positive semidefiniteness; actually, any particular norm and any nonzero norm-constant will work.

<sup>A.8</sup>The symmetrization of  $A$  is  $(A + A^T)/2$ .  $\lambda(\frac{1}{2}(A + A^T)) = \lambda(A + A^T)/2$ .

- By positive semidefiniteness of  $A \in \mathbb{R}^{n \times n}$  we mean,<sup>A.9</sup> [303, §1.3.1] (confer §A.3.1.0.1)

$$x^T A x \geq 0 \quad \forall x \in \mathbb{R}^n \Leftrightarrow A + A^T \succeq 0 \Leftrightarrow \lambda(A + A^T) \succeq 0 \quad (1598)$$

- (§2.9.0.1)

$$A \succeq 0 \Rightarrow A^T = A \quad (1599)$$

$$A \succeq B \Leftrightarrow A - B \succeq 0 \Leftrightarrow A \succeq 0 \text{ or } B \succeq 0 \quad (1600)$$

$$x^T A x \geq 0 \quad \forall x \Leftrightarrow A^T = A \quad (1601)$$

- Matrix symmetry is not intrinsic to positive semidefiniteness;

$$A^T = A, \quad \lambda(A) \succeq 0 \Rightarrow x^T A x \geq 0 \quad \forall x \quad (1602)$$

$$\lambda(A) \succeq 0 \Leftrightarrow A^T = A, \quad x^T A x \geq 0 \quad \forall x \quad (1603)$$

- If  $A^T = A$  then

$$\begin{aligned} \lambda(A) \succeq 0 &\Leftrightarrow A \succeq 0 \\ \lambda(A) \succ 0 &\Leftrightarrow A \succ 0 \end{aligned} \quad (1604)$$

meaning, matrix  $A$  belongs to the positive semidefinite cone (interior) in the subspace of symmetric matrices if and only if its eigenvalues belong to the nonnegative orthant (interior).

$$\langle A, A \rangle = \langle \lambda(A), \lambda(A) \rangle \quad (45)$$

- For  $\mu \in \mathbb{R}$ ,  $A \in \mathbb{R}^{n \times n}$ , and vector  $\lambda(A) \in \mathbb{C}^n$  holding the ordered eigenvalues of  $A$

$$\lambda(\mu I + A) = \mu \mathbf{1} + \lambda(A) \quad (1605)$$

**Proof.**  $A = M J M^{-1}$  and  $\mu I + A = M(\mu I + J)M^{-1}$  where  $J$  is the *Jordan form* for  $A$ ; [368, §5.6, App.B] *id est*,  $\delta(J) = \lambda(A)$ , so  $\lambda(\mu I + A) = \delta(\mu I + J)$  because  $\mu I + J$  is also a Jordan form. ♦

By similar reasoning,

$$\lambda(I + \mu A) = \mathbf{1} + \lambda(\mu A) \quad (1606)$$

For vector  $\sigma(A)$  holding the singular values of any matrix  $A$

$$\sigma(I + \mu A^T A) = \pi(|\mathbf{1} + \mu \sigma(A^T A)|) \quad (1607)$$

$$\sigma(\mu I + A^T A) = \pi(|\mu \mathbf{1} + \sigma(A^T A)|) \quad (1608)$$

where  $\pi$  is the nonlinear permutation-operator sorting its vector argument into nonincreasing order.

- For  $A \in \mathbb{S}^M$  and each and every  $\|w\| = 1$  [228, §7.7 prob.9]

$$w^T A w \leq \mu \Leftrightarrow A \preceq \mu I \Leftrightarrow \lambda(A) \preceq \mu \mathbf{1} \quad (1609)$$

- [228, §2.5.4] (confer (44))

$$A \text{ is normal matrix} \Leftrightarrow \|A\|_F^2 = \lambda(A)^T \lambda(A) \quad (1610)$$

---

<sup>A.9</sup>Strang agrees [368, p.334] it is not  $\lambda(A)$  that requires observation. Yet he is mistaken by proposing the Hermitian part alone  $x^H(A + A^H)x$  be tested, because the anti-Hermitian part does not vanish under complex test unless  $A$  is Hermitian. (1589)

- For  $A \in \mathbb{R}^{m \times n}$

$$A^T A \succeq 0, \quad AA^T \succeq 0 \quad (1611)$$

because, for dimensionally compatible vector  $x$ ,  
 $x^T A^T A x = \|Ax\|_2^2, \quad x^T A A^T x = \|A^T x\|_2^2$ .

- For  $A \in \mathbb{R}^{n \times n}$  and  $c \in \mathbb{R}$

$$\text{tr}(cA) = c \text{ tr}(A) = c \mathbf{1}^T \lambda(A) \quad (\text{\S A.1.1 no.4})$$

For  $m$  a nonnegative integer, (2044)

$$\det(A^m) = \prod_{i=1}^n \lambda(A)_i^m \quad (1612)$$

$$\text{tr}(A^m) = \sum_{i=1}^n \lambda(A)_i^m \quad (1613)$$

- For  $A$  diagonalizable (§A.5),  $A = S\Lambda S^{-1}$ , (confer [368, p.255])

$$\text{rank } A = \text{rank } \delta(\lambda(A)) = \text{rank } \Lambda \quad (1614)$$

meaning, rank is equal to the number of nonzero eigenvalues in vector

$$\lambda(A) \triangleq \delta(\Lambda) \quad (1615)$$

by the 0 eigenvalues theorem (§A.7.3.0.1).

- (Ky Fan) For  $A, B \in \mathbb{S}^n$  [58, §1.2] (confer (1913))

$$\text{tr}(AB) \leq \lambda(A)^T \lambda(B) \quad (1899)$$

with equality (Theobald) when  $A$  and  $B$  are simultaneously diagonalizable [228] with the same ordering of eigenvalues.

- For  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times m}$

$$\text{tr}(AB) = \text{tr}(BA) \quad (1616)$$

and  $\eta$  eigenvalues of the product and commuted product are identical, including their multiplicity; [228, §1.3.20] *id est*,

$$\lambda(AB)_{1:\eta} = \lambda(BA)_{1:\eta}, \quad \eta \triangleq \min\{m, n\} \quad (1617)$$

Any eigenvalues remaining are zero. By the 0 eigenvalues theorem (§A.7.3.0.1),

$$\text{rank}(AB) = \text{rank}(BA), \quad AB \text{ and } BA \text{ diagonalizable} \quad (1618)$$

- For any dimensionally compatible matrices  $A, B$  [228, §0.4]

$$\min\{\text{rank } A, \text{rank } B\} \geq \text{rank}(AB) \quad (1619)$$

- For  $A, B \in \mathbb{S}_+^n$  (confer (262))

$$\text{rank } A + \text{rank } B \geq \text{rank}(A + B) \geq \min\{\text{rank } A, \text{rank } B\} \geq \text{rank}(AB) \quad (1620)$$

- For linearly independent matrices  $A, B \in \mathbb{S}_+^n$   
 $(\S 2.1.2, \mathcal{R}(A) \cap \mathcal{R}(B) = \mathbf{0}, \mathcal{R}(A^T) \cap \mathcal{R}(B^T) = \mathbf{0}, \S B.1.1),$

$$\text{rank } A + \text{rank } B = \text{rank}(A + B) > \min\{\text{rank } A, \text{rank } B\} \geq \text{rank}(AB) \quad (1621)$$

- Because  $\mathcal{R}(A^T A) = \mathcal{R}(A^T)$  and  $\mathcal{R}(A A^T) = \mathcal{R}(A)$  (p.515), for any  $A \in \mathbb{R}^{m \times n}$

$$\text{rank}(A A^T) = \text{rank}(A^T A) = \text{rank } A = \text{rank } A^T \quad (1622)$$

- For  $A \in \mathbb{R}^{m \times n}$  having no nullspace, and for any  $B \in \mathbb{R}^{n \times k}$

$$\text{rank}(AB) = \text{rank}(B) \quad (1623)$$

**Proof.** For any dimensionally compatible matrix  $C$ ,  $\mathcal{N}(CAB) \supseteq \mathcal{N}(AB) \supseteq \mathcal{N}(B)$  is obvious. By assumption  $\exists A^\dagger \ni A^\dagger A = I$ . Let  $C = A^\dagger$ , then  $\mathcal{N}(AB) = \mathcal{N}(B)$  and the stated result follows by conservation of dimension (1741). ♦

- For  $A \in \mathbb{S}^n$  and any nonsingular matrix  $Y$

$$\text{inertia}(A) = \text{inertia}(YAY^T) \quad (1624)$$

a.k.a *Sylvester's law of inertia* (1668) [126, §2.4.3] or *congruence transformation*.

- For  $A, B \in \mathbb{R}^{n \times n}$  square, [228, §0.3.5]

$$\det(AB) = \det(BA) \quad (1625)$$

$$\det(AB) = \det A \det B \quad (1626)$$

Yet for  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times m}$  [84, p.72]

$$\det(I + AB) = \det(I + BA) \quad (1627)$$

- For  $A, B \in \mathbb{S}^n$ , product  $AB$  is symmetric iff  $AB$  is commutative;

$$(AB)^T = AB \Leftrightarrow AB = BA \quad (1628)$$

**Proof.** ( $\Rightarrow$ ) Suppose  $AB = (AB)^T$ .

$$(AB)^T = B^T A^T = BA. \quad AB = (AB)^T \Rightarrow AB = BA.$$

( $\Leftarrow$ ) Suppose  $AB = BA$ .

$$BA = B^T A^T = (AB)^T. \quad AB = BA \Rightarrow AB = (AB)^T. \quad \blacklozenge$$

Matrix symmetry alone is insufficient for product symmetry. Commutativity alone is insufficient for product symmetry. [368, p.26]

- Diagonalizable matrices  $A, B \in \mathbb{R}^{n \times n}$  commute if and only if they are simultaneously diagonalizable. [228, §1.3.12] A product of diagonal matrices is always commutative.

- For  $A, B \in \mathbb{R}^{n \times n}$  and  $AB = BA$

$$x^T A x \geq 0, x^T B x \geq 0 \quad \forall x \Rightarrow \lambda(A + A^T)_i \lambda(B + B^T)_i \geq 0 \quad \forall i \nLeftrightarrow x^T A B x \geq 0 \quad \forall x \quad (1629)$$

the negative result arising because of the schism between the product of eigenvalues  $\lambda(A + A^T)_i \lambda(B + B^T)_i$  and the eigenvalues of the symmetrized matrix product  $\lambda(AB + (AB)^T)_i$ . For example,  $X^2$  is generally not positive semidefinite unless matrix  $X$  is symmetric; then (1611) applies. Simply substituting symmetric matrices changes the outcome:

- For  $A, B \in \mathbb{S}^n$  and  $AB = BA$

$$A \succeq 0, B \succeq 0 \Rightarrow \lambda(AB)_i = \lambda(A)_i \lambda(B)_i \geq 0 \quad \forall i \Leftrightarrow AB \succeq 0 \quad (1630)$$

Positive semidefiniteness of commutative  $A$  and  $B$  is sufficient but not necessary for positive semidefiniteness of product  $AB$ .

**Proof.** Because all symmetric matrices are diagonalizable, (§A.5.1) [368, §5.6] we have  $A = S\Lambda S^T$  and  $B = T\Delta T^T$ , where  $\Lambda$  and  $\Delta$  are real diagonal matrices while  $S$  and  $T$  are orthogonal matrices. Because  $(AB)^T = AB$ , then  $T$  must equal  $S$ , [228, §1.3] and the eigenvalues of  $A$  are ordered in the same way as those of  $B$ ; *id est*,  $\lambda(A)_i = \delta(\Lambda)_i$  and  $\lambda(B)_i = \delta(\Delta)_i$  correspond to the same eigenvector.

$(\Rightarrow)$  Assume  $\lambda(A)_i \lambda(B)_i \geq 0$  for  $i = 1 \dots n$ .  $AB = S\Lambda\Delta S^T$  is symmetric and has nonnegative eigenvalues contained in diagonal matrix  $\Lambda\Delta$  by assumption; hence positive semidefinite by (1598). Now assume  $A, B \succeq 0$ . That, of course, implies  $\lambda(A)_i \lambda(B)_i \geq 0$  for all  $i$  because all the individual eigenvalues are nonnegative.

$(\Leftarrow)$  Suppose  $AB = S\Lambda\Delta S^T \succeq 0$ . Then  $\Lambda\Delta \succeq 0$  by (1598), and so all products  $\lambda(A)_i \lambda(B)_i$  must be nonnegative; meaning,  $\text{sgn}(\lambda(A)) = \text{sgn}(\lambda(B))$ . We may, therefore, conclude nothing about the semidefiniteness of  $A$  and  $B$ .  $\blacklozenge$

- For  $A, B \in \mathbb{S}^n$  and  $A \succeq 0, B \succeq 0$  (Example A.2.1.0.1)

$$AB = BA \Rightarrow \lambda(AB)_i = \lambda(A)_i \lambda(B)_i \geq 0 \quad \forall i \Rightarrow AB \succeq 0 \quad (1631)$$

$$AB = BA \Rightarrow \lambda(AB)_i \geq 0, \lambda(A)_i \lambda(B)_i \geq 0 \quad \forall i \Leftrightarrow AB \succeq 0 \quad (1632)$$

- For  $A, B \in \mathbb{S}^n$  [228, §7.7 prob.3] [229, §4.2.13, §5.2.1]

$$A \succeq 0, B \succeq 0 \Rightarrow A \otimes B \succeq 0 \quad (1633)$$

$$A \succeq 0, B \succeq 0 \Rightarrow A \circ B \succeq 0 \quad (1634)$$

$$A \succ 0, B \succ 0 \Rightarrow A \otimes B \succ 0 \quad (1635)$$

$$A \succ 0, B \succ 0 \Rightarrow A \circ B \succ 0 \quad (1636)$$

where Kronecker and Hadamard products are symmetric.

- For  $A, B \in \mathbb{S}^n$ , (1604)  $A \succeq 0 \Leftrightarrow \lambda(A) \succeq 0$  yet

$$A \succeq 0 \Rightarrow \delta(A) \succeq 0 \quad (1637)$$

$$A \succeq 0 \Rightarrow \text{tr } A \geq 0 \quad (1638)$$

$$A \succeq 0, B \succeq 0 \Rightarrow \text{tr } A \text{ tr } B \geq \text{tr}(AB) \geq 0 \quad (1639)$$

[455, §6.2] Because  $A \succeq 0, B \succeq 0 \Rightarrow \lambda(AB) = \lambda(\sqrt{A}B\sqrt{A}) \succeq 0$  by (1617) and Corollary A.3.1.0.5, then we have  $\text{tr}(AB) \geq 0$ .

$$A \succeq 0 \Leftrightarrow \text{tr}(AB) \geq 0 \quad \forall B \succeq 0 \quad (384)$$

- For  $A, B, C \in \mathbb{S}^n$  (Löwner)

$$\begin{aligned} A \preceq B, B \preceq C &\Rightarrow A \preceq C && \text{(transitivity)} \\ A \preceq B &\Leftrightarrow A + C \preceq B + C && \text{(additivity)} \\ A \preceq B, A \succeq B &\Rightarrow A = B && \text{(antisymmetry)} \\ A \preceq A &&& \text{(reflexivity)} \end{aligned} \quad (1640)$$

$$\begin{aligned} A \preceq B, B \prec C &\Rightarrow A \prec C && \text{(strict transitivity)} \\ A \prec B &\Leftrightarrow A + C \prec B + C && \text{(strict additivity)} \end{aligned} \quad (1641)$$

- For  $A, B \in \mathbb{R}^{n \times n}$

$$x^T A x \geq x^T B x \quad \forall x \Rightarrow \text{tr } A \geq \text{tr } B \quad (1642)$$

**Proof.**  $x^T A x \geq x^T B x \quad \forall x \Leftrightarrow \lambda((A - B) + (A - B)^T)/2 \succeq 0 \Rightarrow \text{tr}(A + A^T - (B + B^T))/2 = \text{tr}(A - B) \geq 0$ . There is no converse.  $\diamond$

- For  $A, B \in \mathbb{S}^n$  [455, §6.2] (Theorem A.3.1.0.4)

$$A \succeq B \Rightarrow \text{tr } A \geq \text{tr } B \quad (1643)$$

$$A \succeq B \Rightarrow \delta(A) \succeq \delta(B) \quad (1644)$$

There is no converse, and restriction to the positive semidefinite cone does not improve the situation. All-strict versions hold.

$$A \succeq B \succeq 0 \Rightarrow \text{rank } A \geq \text{rank } B \quad (1645)$$

$$A \succeq B \succeq 0 \Rightarrow \det A \geq \det B \geq 0 \quad (1646)$$

$$A \succ B \succeq 0 \Rightarrow \det A > \det B \geq 0 \quad (1647)$$

- For  $A, B \in \text{intr } \mathbb{S}_+^n$  [35, §4.2] [228, §7.7.4]

$$A \succeq B \Leftrightarrow A^{-1} \preceq B^{-1}, \quad A \succ 0 \Leftrightarrow A^{-1} \succ 0 \quad (1648)$$

- For  $A, B \in \mathbb{S}^n$  [455, §6.2]

$$\begin{aligned} A \succeq B \succeq 0 &\Rightarrow \sqrt{A} \succeq \sqrt{B} \\ A \succeq 0 &\Leftarrow A^{1/2} \succeq 0 \end{aligned} \quad (1649)$$

- For  $A, B \in \mathbb{S}^n$  and  $AB = BA$  [455, §6.2 prob.3]

$$A \succeq B \succeq 0 \Rightarrow A^k \succeq B^k, \quad k=1, 2, \dots \quad (1650)$$

### A.3.1.0.1 Theorem. Positive semidefinite ordering of eigenvalues.

For  $A, B \in \mathbb{R}^{M \times M}$ , place the eigenvalues of each symmetrized matrix into the respective vectors  $\lambda\left(\frac{1}{2}(A + A^T)\right), \lambda\left(\frac{1}{2}(B + B^T)\right) \in \mathbb{R}^M$ . Then [368, §6]

$$x^T A x \geq 0 \quad \forall x \Leftrightarrow \lambda(A + A^T) \succeq 0 \quad (1651)$$

$$x^T A x > 0 \quad \forall x \neq 0 \Leftrightarrow \lambda(A + A^T) \succ 0 \quad (1652)$$

because  $x^T(A - A^T)x = 0$ . (1586)

Now arrange entries of  $\lambda\left(\frac{1}{2}(A + A^T)\right)$  and  $\lambda\left(\frac{1}{2}(B + B^T)\right)$  in nonincreasing order so  $\lambda\left(\frac{1}{2}(A + A^T)\right)_1$  holds the largest eigenvalue of symmetrized  $A$  while  $\lambda\left(\frac{1}{2}(B + B^T)\right)_1$  holds the largest eigenvalue of symmetrized  $B$ , and so on. Then [228, §7.7 prob.1 prob.9] for  $\kappa \in \mathbb{R}$

$$\begin{aligned} x^T A x \geq x^T B x \quad \forall x &\Rightarrow \lambda(A + A^T) \succeq \lambda(B + B^T) \\ x^T A x \geq x^T I x \kappa \quad \forall x &\Leftrightarrow \lambda\left(\frac{1}{2}(A + A^T)\right) \succeq \kappa \mathbf{1} \end{aligned} \quad (1653)$$

Now let  $A, B \in \mathbb{S}^M$  have diagonalizations  $A = Q\Lambda Q^T$  and  $B = U\Upsilon U^T$  with  $\lambda(A) = \delta(\Lambda)$  and  $\lambda(B) = \delta(\Upsilon)$  arranged in nonincreasing order. Then

$$\begin{aligned} A \succeq B &\Leftrightarrow \lambda(A - B) \succeq 0 \\ A \succeq B &\Rightarrow \lambda(A) \succeq \lambda(B) \\ A \succeq B &\Leftrightarrow \lambda(A) \succeq \lambda(B) \\ S^T A S \succeq B &\Leftrightarrow \lambda(A) \succeq \lambda(B) \end{aligned} \quad (1654)$$

where  $S = Q U^T$ . [455, §7.5] A.10

$\diamond$

---

A.10 ( $\Rightarrow$ )  $S^T A S \succeq B \Rightarrow \lambda(S^T A S) \succeq \lambda(B)$ . But  $S^T A S$  is a matrix similar to  $A$ ; meaning  $\lambda(S^T A S) = \lambda(A)$ .

**A.3.1.0.2 Theorem.** (Weyl) *Eigenvalues of sum.* [228, §4.3.1]

For  $A, B \in \mathbb{R}^{M \times M}$ , place the eigenvalues of each symmetrized matrix into the respective vectors  $\lambda\left(\frac{1}{2}(A + A^T)\right), \lambda\left(\frac{1}{2}(B + B^T)\right) \in \mathbb{R}^M$  in nonincreasing order so  $\lambda\left(\frac{1}{2}(A + A^T)\right)_1$  holds the largest eigenvalue of symmetrized  $A$  while  $\lambda\left(\frac{1}{2}(B + B^T)\right)_1$  holds the largest eigenvalue of symmetrized  $B$ , and so on. Then, for any  $k \in \{1 \dots M\}$

$$\lambda(A + A^T)_k + \lambda(B + B^T)_M \leq \lambda((A + A^T) + (B + B^T))_k \leq \lambda(A + A^T)_k + \lambda(B + B^T)_1 \quad (1655)$$

◊

Weyl's theorem establishes: concavity of the smallest eigenvalue  $\lambda_M$  of a symmetric matrix and convexity of the largest  $\lambda_1$ , via (502), and positive semidefiniteness of a sum of positive semidefinite matrices; for  $A, B \in \mathbb{S}_+^M$

$$\lambda(A)_k + \lambda(B)_M \leq \lambda(A + B)_k \leq \lambda(A)_k + \lambda(B)_1 \quad (1656)$$

Because  $\mathbb{S}_+^M$  is a convex cone (§2.9.0.0.1), then by (178)

$$A, B \succeq 0 \Rightarrow \zeta A + \xi B \succeq 0 \text{ for all } \zeta, \xi \geq 0 \quad (1657)$$

**A.3.1.0.3 Corollary.** *Eigenvalues of sum and difference.* [228, §4.3]

For  $A \in \mathbb{S}^M$  and  $B \in \mathbb{S}_+^M$ , place the eigenvalues of each matrix into respective vectors  $\lambda(A), \lambda(B) \in \mathbb{R}^M$  in nonincreasing order so  $\lambda(A)_1$  holds the largest eigenvalue of  $A$  while  $\lambda(B)_1$  holds the largest eigenvalue of  $B$ , and so on. Then, for any  $k \in \{1 \dots M\}$

$$\lambda(A - B)_k \leq \lambda(A)_k \leq \lambda(A + B)_k \quad (1658)$$

◊

When  $B$  is rank-one positive semidefinite, the eigenvalues interlace; *id est*, for  $B = qq^T$

$$\lambda(A)_{k-1} \leq \lambda(A - qq^T)_k \leq \lambda(A)_k \leq \lambda(A + qq^T)_k \leq \lambda(A)_{k+1} \quad (1659)$$

**A.3.1.0.4 Theorem.** *Positive (semi)definite principal submatrices.* [A.11](#)

- $A \in \mathbb{S}^M$  is positive semidefinite if and only if all  $M$  principal submatrices of dimension  $M-1$  are positive semidefinite and  $\det A$  is nonnegative.
- $A \in \mathbb{S}^M$  is positive definite if and only if any one principal submatrix of dimension  $M-1$  is positive definite and  $\det A$  is positive. ◊

If any one principal submatrix of dimension  $M-1$  is not positive definite, conversely, then  $A$  can neither be. Regardless of symmetry, if  $A \in \mathbb{R}^{M \times M}$  is positive (semi)definite, then the determinant of each and every principal submatrix is (nonnegative) positive. [303, §1.3.1]

---

[A.11](#) A recursive condition for positive (semi)definiteness, this theorem is a synthesis of facts from [228, §7.2] [368, §6.3] (*confer* [303, §1.3.1]).

**A.3.1.0.5 Corollary.** *Positive (semi)definite symmetric products.* [228, p.399]

- (a) If  $A \in \mathbb{S}^M$  is positive definite and any particular dimensionally compatible matrix  $Z$  has no nullspace, then  $Z^T A Z$  is positive definite.
- (b) If matrix  $A \in \mathbb{S}^M$  is positive (semi)definite then, for any matrix  $Z$  of compatible dimension,  $Z^T A Z$  is positive semidefinite.
- (c)  $A \in \mathbb{S}^M$  is positive (semi)definite if and only if there exists a nonsingular  $Z$  such that  $Z^T A Z$  is positive (semi)definite.
- (d) If  $A \in \mathbb{S}^M$  is positive semidefinite and singular, then it remains possible that  $Z^T A Z$  becomes positive definite for some thin  $Z \in \mathbb{R}^{M \times N}$  ( $N < M$ ). [A.12](#)  $\diamond$

We can deduce from these, given nonsingular matrix  $Z$  and any particular dimensionally compatible  $Y$ : matrix  $A \in \mathbb{S}^M$  is positive semidefinite if and only if  $\begin{bmatrix} Z^T \\ Y^T \end{bmatrix} A \begin{bmatrix} Z & Y \end{bmatrix}$  is positive semidefinite. In other words, from the Corollary it follows: for dimensionally compatible  $Z$

$$A \succeq 0 \Leftrightarrow Z^T A Z \succeq 0 \text{ and } Z^T \text{ has a left inverse} \quad (1660)$$

Products such as  $Z^\dagger Z$  and  $ZZ^\dagger$  are symmetric and positive semidefinite although, given  $A \succeq 0$ ,  $Z^\dagger A Z$  and  $Z A Z^\dagger$  are neither necessarily symmetric or positive semidefinite.

**A.3.1.0.6 Theorem.** *Symmetric projector semidefinite.* [21, §III] [22, §6] [249, p.55]  
For symmetric idempotent matrices  $P$  and  $R$

$$\begin{aligned} P, R \succeq 0 \\ P \succeq R \Leftrightarrow \mathcal{R}(P) \supseteq \mathcal{R}(R) \Leftrightarrow \mathcal{N}(P) \subseteq \mathcal{N}(R) \end{aligned} \quad (1661)$$

Projector  $P$  is never positive definite [370, §6.5 prob.20] unless it is the Identity matrix.  $\diamond$

**A.3.1.0.7 Theorem.** *Symmetric positive semidefinite.* [228, p.400]  
Given real matrix  $\Psi$  with rank  $\Psi = 1$

$$\Psi \succeq 0 \Leftrightarrow \Psi = uu^T \quad (1662)$$

is a symmetric dyad where  $u$  is some real vector; *id est*, symmetry is necessary and sufficient for positive semidefiniteness of a rank-1 matrix.  $\diamond$

**Proof.** Any rank-one matrix must have the form  $\Psi = uv^T$ . (§B.1) Suppose  $\Psi$  is symmetric; *id est*,  $v = u$ . For all  $y \in \mathbb{R}^M$ ,  $y^T u u^T y \geq 0$ . Conversely, suppose  $uv^T$  is positive semidefinite. We know that can hold if and only if  $uv^T + vu^T \succeq 0 \Leftrightarrow$  for all normalized  $y \in \mathbb{R}^M$ ,  $2y^T u v^T y \geq 0$ ; but that is possible only if  $v = u$ .  $\blacklozenge$

The same does not hold true for matrices of higher rank, as Example A.2.1.0.1 shows.

---

[A.12](#)This means coefficients, of orthogonal projection of vectorized  $A$  on a subset of extreme directions from  $\mathbb{S}_+^M$  determined by  $Z$ , can be positive (by the interpretation in §E.6.4.3).

## A.4 Schur complement

Consider *Schur-form* partitioned matrix  $G$ : Given  $A^T = A$  and  $C^T = C$ , then [63]

$$\begin{aligned} G &= \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0 & \text{(a)} \\ \Leftrightarrow A &\succeq 0, \quad B^T(I - AA^\dagger) = \mathbf{0}, \quad C - B^T A^\dagger B \succeq 0 & \text{(b)} \\ \Leftrightarrow C &\succeq 0, \quad B(I - CC^\dagger) = \mathbf{0}, \quad A - BC^\dagger B^T \succeq 0 & \text{(c)} \end{aligned} \quad (1663)$$

where  $A^\dagger$  denotes the Moore-Penrose (pseudo)inverse (§E). In the first instance,  $I - AA^\dagger$  is a symmetric projection matrix orthogonally projecting on  $\mathcal{N}(A^T)$ . (2094) It is apparently required

$$\mathcal{R}(B) \perp \mathcal{N}(A^T) \quad (1664)$$

which precludes  $A = \mathbf{0}$  when  $B$  is any nonzero matrix. Note that  $A \succ 0 \Rightarrow A^\dagger = A^{-1}$ ; thereby, the projection matrix vanishes. Likewise, in the second instance,  $I - CC^\dagger$  projects orthogonally on  $\mathcal{N}(C^T)$ . It is required

$$\mathcal{R}(B^T) \perp \mathcal{N}(C^T) \quad (1665)$$

which precludes  $C = \mathbf{0}$  for  $B$  nonzero. Again,  $C \succ 0 \Rightarrow C^\dagger = C^{-1}$ . So we get, for  $A$  or  $C$  nonsingular,

$$\begin{aligned} G &= \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0 \\ \Leftrightarrow A &\succ 0, \quad C - B^T A^{-1} B \succeq 0 \\ \Leftrightarrow C &\succ 0, \quad A - B C^{-1} B^T \succeq 0 \end{aligned} \quad (1666)$$

When  $A$  is full-rank then, for all  $B$  of compatible dimension,  $\mathcal{R}(B)$  is in  $\mathcal{R}(A)$ . Likewise, when  $C$  is full-rank,  $\mathcal{R}(B^T)$  is in  $\mathcal{R}(C)$ . Thus, for  $A$  and  $C$  nonsingular,

$$\begin{aligned} G &= \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succ 0 \\ \Leftrightarrow A &\succ 0, \quad C - B^T A^{-1} B \succ 0 \\ \Leftrightarrow C &\succ 0, \quad A - B C^{-1} B^T \succ 0 \end{aligned} \quad (1667)$$

where  $C - B^T A^{-1} B$  is called the *Schur complement of  $A$  in  $G$* , while the *Schur complement of  $C$  in  $G$*  is  $A - B C^{-1} B^T$ . [174, §4.8]

Origin of the term *Schur complement* is from complementary *inertia*: [126, §2.4.4] Define

$$\text{inertia}\left(G \in \mathbb{S}^M\right) \triangleq \{p, z, n\} \quad (1668)$$

where  $p, z, n$  respectively represent number of positive, zero, and negative eigenvalues of  $G$ ; *id est*,

$$M = p + z + n \quad (1669)$$

Then, when  $A$  is invertible,

$$\text{inertia}(G) = \text{inertia}(A) + \text{inertia}(C - B^T A^{-1} B) \quad (1670)$$

and when  $C$  is invertible,

$$\text{inertia}(G) = \text{inertia}(C) + \text{inertia}(A - B C^{-1} B^T) \quad (1671)$$

**A.4.0.0.1 Example.** *Equipartition inertia.*

[58, §1.2 exer.17]

When  $A = C = \mathbf{0}$ , denoting nonincreasingly ordered singular values of matrix  $B \in \mathbb{R}^{m \times m}$  by  $\sigma(B) \in \mathbb{R}_+^m$ , then we have eigenvalues

$$\lambda(G) = \lambda\left(\begin{bmatrix} \mathbf{0} & B \\ B^T & \mathbf{0} \end{bmatrix}\right) = \begin{bmatrix} \sigma(B) \\ -\Xi\sigma(B) \end{bmatrix} \quad (1672)$$

and

$$\text{inertia}(G) = \text{inertia}(B^T B) + \text{inertia}(-B^T B) \quad (1673)$$

where  $\Xi$  is the order-reversing permutation matrix defined in (1900).  $\square$ **A.4.0.0.2 Example.** *Nonnegative polynomial.*

[35, p.163]

Quadratic multivariate polynomial  $x^T A x + 2b^T x + c$  is a convex function of vector  $x$  if and only if  $A \succeq 0$ , but sublevel set  $\{x \mid x^T A x + 2b^T x + c \leq 0\}$  is convex if  $A \succeq 0$  yet not *vice versa*. Schur-form positive semidefiniteness is sufficient for polynomial convexity but necessary and sufficient for nonnegativity:

$$\begin{bmatrix} A & b \\ b^T & c \end{bmatrix} \succeq 0 \Leftrightarrow [x^T \ 1] \begin{bmatrix} A & b \\ b^T & c \end{bmatrix} \begin{bmatrix} x \\ 1 \end{bmatrix} \geq 0 \Leftrightarrow x^T A x + 2b^T x + c \geq 0 \quad \forall x \quad (1674)$$

Everything here is extensible to univariate polynomials; e.g.,  $x \triangleq [t^n \ t^{n-1} \ t^{n-2} \cdots \ t]^T$ . $\square$ **A.4.0.0.3 Example.** *Schur-form fractional function trace minimization.*

From (1638),

$$\begin{aligned} \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0 &\Rightarrow \text{tr}(A + C) \geq 0 \\ \Updownarrow \\ \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0}^T & C - B^T A^{-1} B \end{bmatrix} \succeq 0 &\Rightarrow \text{tr}(C - B^T A^{-1} B) \geq 0 \\ \Updownarrow \\ \begin{bmatrix} A - B C^{-1} B^T & \mathbf{0} \\ \mathbf{0}^T & C \end{bmatrix} \succeq 0 &\Rightarrow \text{tr}(A - B C^{-1} B^T) \geq 0 \end{aligned} \quad (1675)$$

Since  $\text{tr}(C - B^T A^{-1} B) \geq 0 \Leftrightarrow \text{tr } C \geq \text{tr}(B^T A^{-1} B) \geq 0$  for example, then minimization of  $\text{tr } C$  is necessary and sufficient for minimization of  $\text{tr}(C - B^T A^{-1} B)$  when both are under constraint  $\begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0$ .  $\square$

**A.4.0.1 Schur-form nullspace basis**

From (1663),

$$\begin{aligned} G = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0 & \\ \Updownarrow \\ \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0}^T & C - B^T A^\dagger B \end{bmatrix} \succeq 0 \quad \text{and} \quad B^T(I - A A^\dagger) = \mathbf{0} & \quad (1676) \\ \Updownarrow \\ \begin{bmatrix} A - B C^\dagger B^T & \mathbf{0} \\ \mathbf{0}^T & C \end{bmatrix} \succeq 0 \quad \text{and} \quad B(I - C C^\dagger) = \mathbf{0} & \end{aligned}$$

These facts plus Moore-Penrose condition (§E.0.1) provide a partial basis:

$$\text{basis } \mathcal{N}\left(\begin{bmatrix} A & B \\ B^T & C \end{bmatrix}\right) \supseteq \begin{bmatrix} I - AA^\dagger & \mathbf{0} \\ \mathbf{0}^T & I - CC^\dagger \end{bmatrix} \quad (1677)$$

#### A.4.0.1.1 Example. Sparse Schur-form.

Setting matrix  $A$  to the Identity simplifies the Schur-form. One consequence relates definiteness of three quantities:

$$\begin{bmatrix} I & B \\ B^T & C \end{bmatrix} \succeq 0 \Leftrightarrow C - B^T B \succeq 0 \Leftrightarrow \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0}^T & C - B^T B \end{bmatrix} \succeq 0 \quad (1678)$$

□

#### A.4.0.1.2 Exercise. Eigenvalues $\lambda$ of sparse Schur-form.

Prove: given  $C - B^T B = \mathbf{0}$ , for  $B \in \mathbb{R}^{m \times n}$  and  $C \in \mathbb{S}^n$

$$\lambda\left(\begin{bmatrix} I & B \\ B^T & C \end{bmatrix}\right)_i = \begin{cases} 1 + \lambda(C)_i, & 1 \leq i \leq n \\ 1, & n < i \leq m \\ 0, & \text{otherwise} \end{cases} \quad (1679)$$

▼

#### A.4.0.1.3 Theorem. Rank of partitioned matrices.

[455, §2.2 prob.7]

When symmetric matrix  $A$  is invertible and  $C$  is symmetric,

$$\begin{aligned} \text{rank}\left[\begin{array}{cc} A & B \\ B^T & C \end{array}\right] &= \text{rank}\left[\begin{array}{cc} A & \mathbf{0} \\ \mathbf{0}^T & C - B^T A^{-1} B \end{array}\right] \\ &= \text{rank } A + \text{rank}(C - B^T A^{-1} B) \end{aligned} \quad (1680)$$

equals rank of main diagonal block  $A$  plus rank of its Schur complement.

Similarly, when symmetric matrix  $C$  is invertible and  $A$  is symmetric,

$$\begin{aligned} \text{rank}\left[\begin{array}{cc} A & B \\ B^T & C \end{array}\right] &= \text{rank}\left[\begin{array}{cc} A - BC^{-1}B^T & \mathbf{0} \\ \mathbf{0}^T & C \end{array}\right] \\ &= \text{rank}(A - BC^{-1}B^T) + \text{rank } C \end{aligned} \quad (1681)$$

◇

**Proof.** The first assertion (1680) holds if and only if [228, §0.4.6c]

$$\exists \text{ nonsingular } X, Y \ni X \left[\begin{array}{cc} A & B \\ B^T & C \end{array}\right] Y = \left[\begin{array}{cc} A & \mathbf{0} \\ \mathbf{0}^T & C - B^T A^{-1} B \end{array}\right] \quad (1682)$$

Let [228, §7.7.6]

$$Y = X^T = \left[\begin{array}{cc} I & -A^{-1}B \\ \mathbf{0}^T & I \end{array}\right] \quad (1683)$$

◆

From Corollary A.3.1.0.3, eigenvalues are related by

$$0 \preceq \lambda(C - B^T A^{-1} B) \preceq \lambda(C) \quad (1684)$$

$$0 \preceq \lambda(A - BC^{-1}B^T) \preceq \lambda(A) \quad (1685)$$

which means

$$\text{rank}(C - B^T A^{-1} B) \leq \text{rank } C \quad (1686)$$

$$\text{rank}(A - B C^{-1} B^T) \leq \text{rank } A \quad (1687)$$

Therefore

$$\text{rank} \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \leq \text{rank } A + \text{rank } C \quad (1688)$$

#### A.4.0.1.4 Lemma. Rank of Schur-form block. [154] [152]

Matrix  $B \in \mathbb{R}^{m \times n}$  has  $\text{rank } B \leq \rho$  if and only if there exist matrices  $A \in \mathbb{S}^m$  and  $C \in \mathbb{S}^n$  such that

$$\text{rank} \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0}^T & C \end{bmatrix} \leq 2\rho \quad \text{and} \quad G = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0 \quad (1689)$$

◊

Schur-form positive semidefiniteness alone implies  $\text{rank } A \geq \text{rank } B$  and  $\text{rank } C \geq \text{rank } B$ . But, even in absence of semidefiniteness, we must always have  $\text{rank } G \geq \text{rank } A, \text{rank } B, \text{rank } C$  by fundamental linear algebra.

### A.4.1 Determinant

$$G = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \quad (1690)$$

We consider again a matrix  $G$  partitioned like (1663), but not necessarily positive (semi)definite, where  $A$  and  $C$  are symmetric.

- When  $A$  is invertible,

$$\det G = \det A \det(C - B^T A^{-1} B) \quad (1691)$$

When  $C$  is invertible,

$$\det G = \det C \det(A - B C^{-1} B^T) \quad (1692)$$

- When  $B$  is full-rank and thin,  $C = \mathbf{0}$ , and  $A \succeq 0$ , then [65, §10.1.1]

$$\det G \neq 0 \Leftrightarrow A + BB^T \succ 0 \quad (1693)$$

When  $B$  is a (column) vector, then for all  $C \in \mathbb{R}$  and all  $A$  of dimension compatible with  $G$

$$\det G = \det(A)C - B^T A_{\text{cof}}^T B \quad (1694)$$

while for  $C \neq 0$

$$\det G = C \det(A - \frac{1}{C}BB^T) \quad (1695)$$

where  $A_{\text{cof}}$  is the matrix of cofactors [368, §4] corresponding to  $A$ .

- When  $B$  is full-rank and wide,  $A = \mathbf{0}$ , and  $C \succeq 0$ , then

$$\det G \neq 0 \Leftrightarrow C + B^T B \succ 0 \quad (1696)$$

When  $B$  is a row vector, then for  $A \neq 0$  and all  $C$  of dimension compatible with  $G$

$$\det G = A \det(C - \frac{1}{A}B^T B) \quad (1697)$$

while for all  $A \in \mathbb{R}$

$$\det G = \det(C)A - BC_{\text{cof}}^T B^T \quad (1698)$$

where  $C_{\text{cof}}$  is the matrix of cofactors corresponding to  $C$ .

## A.5 Eigenvalue decomposition

All square matrices  $A$  have associated eigenvalues  $\lambda$  and eigenvectors  $S$ ; [A.13](#) if not square,  $As_i = \lambda_i s_i$  becomes impossible dimensionally. Eigenvectors must be nonzero.

When a square matrix  $X \in \mathbb{R}^{m \times m}$  is *diagonalizable*, [368, §5.6] then

$$X = S\Lambda S^{-1} = [s_1 \cdots s_m] \Lambda \begin{bmatrix} w_1^T \\ \vdots \\ w_m^T \end{bmatrix} = \sum_{i=1}^m \lambda_i s_i w_i^T \quad (1699)$$

where  $\{s_i \in \mathcal{N}(X - \lambda_i I) \subseteq \mathbb{C}^m\}$  are l.i. (right-)eigenvectors constituting the columns of  $S \in \mathbb{C}^{m \times m}$  defined by

$$XS = S\Lambda \quad \text{rather} \quad Xs_i \triangleq \lambda_i s_i, \quad i = 1 \dots m \quad (1700)$$

$\{w_i \in \mathcal{N}(X^T - \lambda_i I) \subseteq \mathbb{C}^m\}$  are linearly independent *left-eigenvectors* of  $X$  (eigenvectors of  $X^T$ ) constituting the rows of  $S^{-1}$  defined by [228]

$$S^{-1}X = \Lambda S^{-1} \quad \text{rather} \quad w_i^T X \triangleq \lambda_i w_i^T, \quad i = 1 \dots m \quad (1701)$$

and where  $\{\lambda_i \in \mathbb{C}\}$  are eigenvalues ([1615](#))

$$\delta(\lambda(X)) = \Lambda \in \mathbb{C}^{m \times m} \quad (1702)$$

corresponding to both left and right eigenvectors; *id est*,  $\lambda(X) = \lambda(X^T)$ .

There is no connection between diagonalizability and invertibility of  $X$ . [368, §5.2] Diagonalizability is guaranteed by a full set of linearly independent eigenvectors, whereas invertibility is guaranteed by all nonzero eigenvalues.

$$\begin{aligned} \text{distinct eigenvalues} &\Rightarrow \text{l.i. eigenvectors} \Leftrightarrow \text{diagonalizable} \\ \text{not diagonalizable} &\Rightarrow \text{repeated eigenvalue} \end{aligned} \quad (1703)$$

$\begin{bmatrix} 1 & 0 & -1 \\ -1 & 1 & 0 \\ 3 & -1 & -2 \end{bmatrix}$  is not diagonalizable, for example, having three 0-eigenvalues which are hard to compute with accuracy better than 1E-6. ([Yates, D'Errico](#))

### A.5.0.0.1 Theorem. Real eigenvector.

Eigenvectors of a real matrix corresponding to real eigenvalues must be real.  $\diamond$

**Proof.**  $Ax = \lambda x$ . Given  $\lambda = \lambda^*$ ,  $x^H Ax = \lambda x^H x = \lambda \|x\|^2 = x^T A x^* \Rightarrow x = x^*$ , where  $x^H = x^{*T}$ . The converse is equally simple.  $\blacklozenge$

### A.5.0.1 Uniqueness

From the *fundamental theorem of algebra*, [387] which guarantees existence of zeros for a given polynomial, it follows: Given a particular square matrix, its eigenvalues and their multiplicity are unique; meaning, there is no other set of eigenvalues for that matrix. (Conversely, many different matrices may share the same unique set of eigenvalues; *e.g.*, for any  $X$ ,  $\lambda(X) = \lambda(X^T)$ .)

Uniqueness of eigenvectors, in contrast, disallows multiplicity of the same direction:

---

[A.13](#)Prefix *eigen* is from the German; in this context meaning, something akin to “characteristic”. [364, p.14]

#### A.5.0.1.1 Definition. Unique eigenvectors.

When eigenvectors are *unique*, we mean: unique to within a real nonzero scaling, and their directions are distinct.  $\triangle$

If  $S$  is a matrix of eigenvectors of  $X$  as in (1699), for example, then  $-S$  is certainly another matrix of eigenvectors decomposing  $X$  with the same eigenvalues. Although directions are distinct, eigenvectors  $-S$  are equivalent to eigenvectors  $S$  by Definition A.5.0.1.1.

For any square matrix, the eigenvector corresponding to a distinct eigenvalue is unique; [364, p.220]

$$\text{distinct eigenvalues} \Rightarrow \text{eigenvectors unique} \quad (1704)$$

Eigenvectors corresponding to a repeated eigenvalue are not unique for a diagonalizable matrix;

$$\text{repeated eigenvalue} \Rightarrow \text{eigenvectors not unique} \quad (1705)$$

Proof follows from the observation: any linear combination of distinct eigenvectors of diagonalizable  $X$ , corresponding to a particular eigenvalue, produces another eigenvector. For eigenvalue  $\lambda$  whose multiplicity [A.14](#)  $\dim \mathcal{N}(X - \lambda I)$  exceeds 1, in other words, any choice of independent vectors from  $\mathcal{N}(X - \lambda I)$  (of the same multiplicity) constitutes eigenvectors corresponding to  $\lambda$ .  $\spadesuit$

*Caveat* diagonalizability insures linear independence which implies existence of distinct eigenvectors. We may conclude, for diagonalizable matrices,

$$\text{distinct eigenvalues} \Leftrightarrow \text{eigenvectors unique} \quad (1706)$$

#### A.5.0.2 Invertible matrix

When diagonalizable matrix  $X \in \mathbb{R}^{m \times m}$  is *nonsingular* (no zero eigenvalues), then it has an inverse obtained simply by inverting eigenvalues in (1699):

$$X^{-1} = S\Lambda^{-1}S^{-1} \quad (1707)$$

#### A.5.0.3 eigenmatrix

The (right-)eigenvectors  $\{s_i\}$  (1699) are naturally orthogonal  $w_i^T s_j = 0$  to left-eigenvectors  $\{w_i\}$  except, for  $i=1 \dots m$ ,  $w_i^T s_i = 1$ ; called a *biorthogonality condition* [406, §2.2.4] [228] because neither set of left or right eigenvectors is necessarily an orthogonal set. Consequently, each dyad from a diagonalization is an independent ([§B.1.1](#)) nonorthogonal projector because

$$s_i w_i^T s_i w_i^T = s_i w_i^T \quad (1708)$$

(whereas the dyads of singular value decomposition are not inherently projectors (*confer*(1716))).

Dyads of eigenvalue decomposition can be termed *eigenmatrices* because

$$X s_i w_i^T = \lambda_i s_i w_i^T \quad (1709)$$

Sum of the eigenmatrices is the Identity;

$$\sum_{i=1}^m s_i w_i^T = I \quad (1710)$$

---

[A.14](#) A matrix is diagonalizable iff *algebraic multiplicity* (number of occurrences of same eigenvalue) equals *geometric multiplicity*  $\dim \mathcal{N}(X - \lambda I) = m - \text{rank}(X - \lambda I)$  [364, p.15] (number of *Jordan blocks* w.r.t  $\lambda$  or number of corresponding l.i. eigenvectors).

### A.5.1 Symmetric matrix diagonalization

The set of *normal matrices* is, precisely, that set of all real matrices having a complete orthonormal set of eigenvectors; [455, §8.1] [370, prob.10.2.31] *id est*, any matrix  $X$  for which  $XX^T = X^TX$ ; [181, §7.1.3] [364, p.3] *e.g.*, symmetric, orthogonal, and circulant matrices [192]. All normal matrices are diagonalizable.

A symmetric matrix is a special normal matrix whose eigenvalues  $\Lambda$  must be real<sup>A.15</sup> and whose eigenvectors  $S$  can be chosen to make a real orthonormal set; [370, §6.4] [368, p.315] *id est*, for  $X \in \mathbb{S}^m$

$$X = S\Lambda S^T = [s_1 \cdots s_m] \Lambda \begin{bmatrix} s_1^T \\ \vdots \\ s_m^T \end{bmatrix} = \sum_{i=1}^m \lambda_i s_i s_i^T \quad (1711)$$

where  $\delta^2(\Lambda) = \Lambda \in \mathbb{S}^m$  (§A.1) and  $S^{-1} = S^T \in \mathbb{R}^{m \times m}$  (orthogonal matrix, §B.5.2) because of symmetry:  $S\Lambda S^{-1} = S^{-T}\Lambda S^T$ . By 0 eigenvalues theorem A.7.3.0.1,

$$\begin{aligned} \mathcal{R}\{s_i \mid \lambda_i \neq 0\} &= \mathcal{R}(A) = \mathcal{R}(A^T) \\ \mathcal{R}\{s_i \mid \lambda_i = 0\} &= \mathcal{N}(A^T) = \mathcal{N}(A) \end{aligned} \quad (1712)$$

#### A.5.1.1 eigenvalue $\lambda$ ordering

Because arrangement of eigenvectors and their corresponding eigenvalues is arbitrary, eigenvalues are often arranged in nonincreasing order; as is the convention for singular value decomposition (§A.6). There are certainly circumstances demanding otherwise; *e.g.*, a direction vector of convex iteration (§4.5.1.1) can require simultaneous diagonalizability.

#### A.5.1.2 diagonal matrix diagonalization

Then to diagonalize a symmetric matrix that is already a diagonal matrix, orthogonal matrix  $S$  becomes a permutation matrix  $\Xi$ . Given vector  $a$ , for example, diagonal matrix  $\delta(a)$  has eigenvalue decomposition

$$\delta(a) = \Xi \delta(\Xi^T a) \Xi^T \quad (1713)$$

#### A.5.1.3 invertible symmetric matrix

When symmetric matrix  $X \in \mathbb{S}^m$  is nonsingular (invertible), then its inverse (obtained by inverting eigenvalues in (1711)) is also symmetric:

$$X^{-1} = S\Lambda^{-1}S^T \in \mathbb{S}^m \quad (1714)$$

#### A.5.1.4 positive semidefinite matrix square root

When  $X \in \mathbb{S}_+^m$ , its unique positive semidefinite matrix square root is defined (1711)

$$\sqrt{X} \triangleq S\sqrt{\Lambda}S^T \in \mathbb{S}_+^m \quad (1715)$$

where the square root of nonnegative diagonal matrix  $\sqrt{\Lambda}$  is taken entrywise and positive. Then  $X = \sqrt{X}\sqrt{X}$ .

---

<sup>A.15</sup>**Proof.** Suppose  $\lambda_i$  is an eigenvalue corresponding to eigenvector  $s_i$  of real  $A = A^T$ . Then  $s_i^H A s_i = s_i^T A s_i^*$  (by transposition)  $\Rightarrow s_i^T \lambda_i s_i = s_i^T \lambda_i^* s_i^*$  because  $(As_i)^* = (\lambda_i s_i)^*$  by assumption. So we have  $\lambda_i \|s_i\|^2 = \lambda_i^* \|s_i\|^2$ . There is no converse. ♦

## A.6 Singular value decomposition, SVD

### A.6.1 Compact SVD

[181, §2.5.4] For any  $A \in \mathbb{R}^{m \times n}$

$$A = U\Sigma Q^T = [u_1 \cdots u_\eta] \Sigma \begin{bmatrix} q_1^T \\ \vdots \\ q_\eta^T \end{bmatrix} = \sum_{i=1}^{\eta} \sigma_i u_i q_i^T \quad (1716)$$

$$U \in \mathbb{R}^{m \times \eta}, \quad \Sigma \in \mathbb{R}_+^{\eta \times \eta}, \quad Q \in \mathbb{R}^{n \times \eta}$$

$$U^T U = I, \quad Q^T Q = I$$

where  $U$  and  $Q$  are always thin-or-square real, each having orthonormal columns, and where

$$\eta \triangleq \min\{m, n\} \quad (1717)$$

Square matrix  $\Sigma$  is diagonal (§A.1.1)

$$\delta^2(\Sigma) = \Sigma \in \mathbb{R}_+^{\eta \times \eta} \quad (1718)$$

holding the singular values  $\{\sigma_i \in \mathbb{R}\}$  of  $A$  which are always arranged in nonincreasing order by convention and are related to eigenvalues  $\lambda$  by A.16

$$\sigma(A)_i = \sigma(A^T)_i = \begin{cases} \sqrt{\lambda(A^T A)_i} = \sqrt{\lambda(A A^T)_i} = \lambda(\sqrt{A^T A})_i = \lambda(\sqrt{A A^T})_i > 0, & 1 \leq i \leq \rho \\ 0, & \rho < i \leq \eta \end{cases} \quad (1719)$$

of which the last  $\eta - \rho$  are 0, A.17 where

$$\rho \triangleq \text{rank } A = \text{rank } \Sigma \quad (1720)$$

A point sometimes lost: Any real matrix may be decomposed in terms of its real singular values  $\sigma(A) \in \mathbb{R}^\eta$  and real matrices  $U$  and  $Q$  as in (1716), where [181, §2.5.3]

$$\begin{aligned} \mathcal{R}\{u_i \mid \sigma_i \neq 0\} &= \mathcal{R}(A) \\ \mathcal{R}\{u_i \mid \sigma_i = 0\} &\subseteq \mathcal{N}(A^T) \\ \mathcal{R}\{q_i \mid \sigma_i \neq 0\} &= \mathcal{R}(A^T) \\ \mathcal{R}\{q_i \mid \sigma_i = 0\} &\subseteq \mathcal{N}(A) \end{aligned} \quad (1721)$$

### A.6.2 Subcompact SVD

Some authors allow only nonzero singular values. In that case the compact decomposition can be made smaller; it can be redimensioned in terms of rank  $\rho$  because, for any  $A \in \mathbb{R}^{m \times n}$

$$\rho = \text{rank } A = \text{rank } \Sigma = \max\{i \in \{1 \dots \eta\} \mid \sigma_i \neq 0\} \leq \eta \quad (1722)$$

- There are  $\eta$  singular values. For any flavor SVD, rank is equivalent to the number of nonzero singular values on the main diagonal of  $\Sigma$ .

**A.16**  $\sigma(A^T A) = \lambda(A^T A)$  and  $\sigma(A A^T) = \lambda(A A^T)$ . (§A.6.3.2) But when matrix  $A$  is normal,  $\sigma(A) = |\lambda(A)|$ . [455, §8.1]

**A.17** For  $\eta = n$ ,  $\sigma(A) = \sqrt{\lambda(A^T A)} = \lambda(\sqrt{A^T A})$ . For  $\eta = m$ ,  $\sigma(A) = \sqrt{\lambda(A A^T)} = \lambda(\sqrt{A A^T})$ .

Now

$$A = U\Sigma Q^T = [u_1 \cdots u_\rho] \Sigma \begin{bmatrix} q_1^T \\ \vdots \\ q_\rho^T \end{bmatrix} = \sum_{i=1}^\rho \sigma_i u_i q_i^T \quad (1723)$$

$$U \in \mathbb{R}^{m \times \rho}, \quad \Sigma \in \mathbb{R}_+^{\rho \times \rho}, \quad Q \in \mathbb{R}^{n \times \rho}$$

$$U^T U = I, \quad Q^T Q = I$$

where the main diagonal of diagonal matrix  $\Sigma$  has no 0 entries, and

$$\begin{aligned} \mathcal{R}\{u_i\} &= \mathcal{R}(A) \\ \mathcal{R}\{q_i\} &= \mathcal{R}(A^T) \end{aligned} \quad (1724)$$

### A.6.3 Full SVD

Another common and useful expression of the SVD makes  $U$  and  $Q$  square; making the decomposition larger than compact SVD. Completing the nullspace bases in  $U$  and  $Q$  from (1721) provides what is called the *full singular value decomposition* of  $A \in \mathbb{R}^{m \times n}$  [368, App.A]. Orthonormal real matrices  $U$  and  $Q$  become orthogonal matrices (§B.5):

$$\begin{aligned} \mathcal{R}\{u_i \mid \sigma_i \neq 0\} &= \mathcal{R}(A) \\ \mathcal{R}\{u_i \mid \sigma_i = 0\} &= \mathcal{N}(A^T) \\ \mathcal{R}\{q_i \mid \sigma_i \neq 0\} &= \mathcal{R}(A^T) \\ \mathcal{R}\{q_i \mid \sigma_i = 0\} &= \mathcal{N}(A) \end{aligned} \quad (1725)$$

For any matrix  $A$  having rank  $\rho$  ( $= \text{rank } \Sigma$ )

$$A = U\Sigma Q^T = [u_1 \cdots u_m] \Sigma \begin{bmatrix} q_1^T \\ \vdots \\ q_n^T \end{bmatrix} = \sum_{i=1}^\eta \sigma_i u_i q_i^T$$

$$= [m \times \rho \text{ basis } \mathcal{R}(A) \quad m \times m-\rho \text{ basis } \mathcal{N}(A^T)] \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \end{bmatrix} \begin{bmatrix} (n \times \rho \text{ basis } \mathcal{R}(A^T))^T \\ (n \times n-\rho \text{ basis } \mathcal{N}(A))^T \end{bmatrix}$$

$$U \in \mathbb{R}^{m \times m}, \quad \Sigma \in \mathbb{R}_+^{m \times n}, \quad Q \in \mathbb{R}^{n \times n}$$

$$U^T = U^{-1}, \quad Q^T = Q^{-1} \quad (1726)$$

where upper limit of summation  $\eta$  is defined in (1717). Matrix  $\Sigma$  is no longer necessarily square, now padded with respect to (1718) by  $m-\eta$  zero rows or  $n-\eta$  zero columns; the nonincreasingly ordered (possibly 0) singular values appear along its main diagonal as for compact SVD (1719).

An important geometrical interpretation of SVD is given in Figure 180 for  $m=n=2$ : The image of the unit sphere under any  $m \times n$  matrix multiplication is an ellipse. Considering the three factors of the SVD separately, note that  $Q^T$  is a pure rotation of the circle. Figure 180 shows how the axes  $q_1$  and  $q_2$  are first rotated by  $Q^T$  to coincide with the coordinate axes. Second, the circle is stretched by  $\Sigma$  in the directions of the coordinate axes to form an ellipse. The third step rotates the ellipse by  $U$  into its final position. Note how  $q_1$  and  $q_2$  are rotated to end up as  $u_1$  and  $u_2$ , the principal axes of the final ellipse. A direct calculation shows that  $Aq_j = \sigma_j u_j$ . Thus  $q_j$  is first rotated to coincide with the  $j^{\text{th}}$  coordinate axis, stretched by a factor  $\sigma_j$ , and then rotated to point in the direction of  $u_j$ . All of this is beautifully illustrated for  $2 \times 2$  matrices by the MATLAB code `eigshow.m` (see [371]).

$$A = U\Sigma Q^T = [u_1 \cdots u_m] \Sigma \begin{bmatrix} q_1^T \\ \vdots \\ q_n^T \end{bmatrix} = \sum_{i=1}^n \sigma_i u_i q_i^T$$

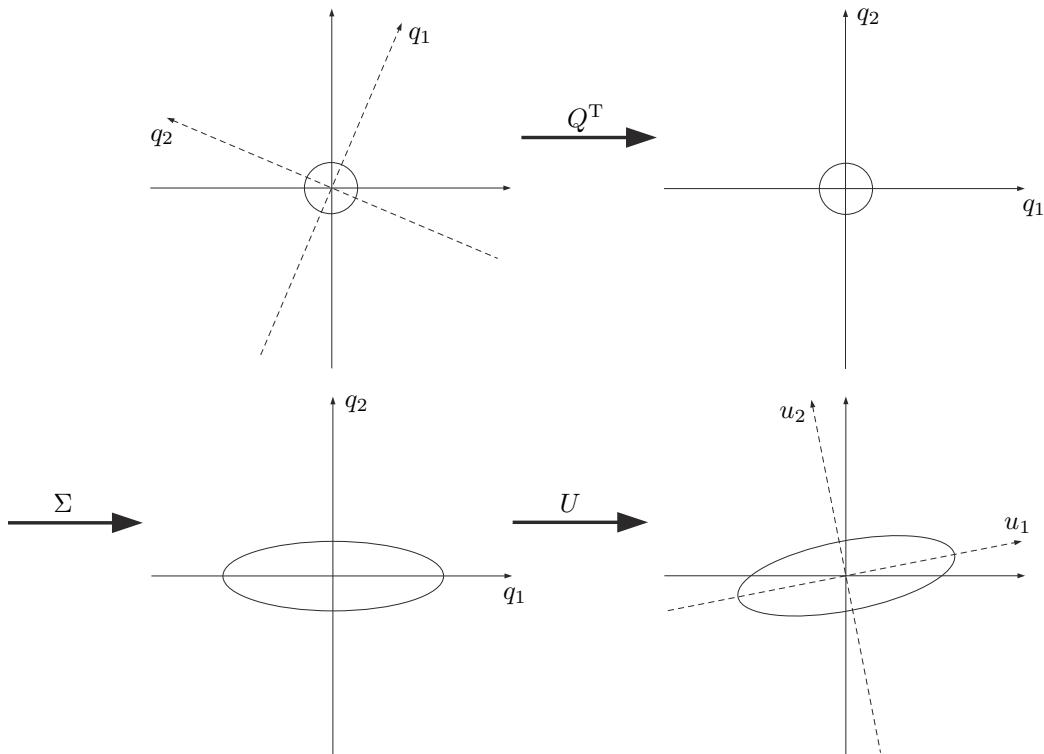


Figure 180: Full SVD geometrical interpretation [302]: Image of circle  $\{x \in \mathbb{R}^2 \mid \|x\|_2 = 1\}$ , under matrix multiplication  $Ax$ , is generally an ellipse. For the example illustrated,  $U \triangleq [u_1 \ u_2] \in \mathbb{R}^{2 \times 2}$ ,  $Q \triangleq [q_1 \ q_2] \in \mathbb{R}^{2 \times 2}$ .

A direct consequence of the geometric interpretation is that the largest singular value  $\sigma_1$  measures the “magnitude” of  $A$  (its 2-norm):

$$\|A\|_2 = \sup_{\|x\|_2=1} \|Ax\|_2 = \sigma_1 \quad (1727)$$

This means that  $\|A\|_2$  is the length of the longest principal semiaxis of the ellipse.

Expressions for  $U$ ,  $Q$ , and  $\Sigma$  follow readily from (1726),

$$AA^T U = U\Sigma\Sigma^T \text{ and } A^T A Q = Q\Sigma^T \Sigma \quad (1728)$$

demonstrating that the columns of  $U$  are the eigenvectors of  $AA^T$  and the columns of  $Q$  are the eigenvectors of  $A^T A$ .  
—Muller, Magaia, & Herbst [302]

### A.6.3.1 SVD of positive semidefinite matrices

From (1719) and (1715) for  $A \succeq 0$

$$\sigma(A)_i = \begin{cases} \sqrt{\lambda(A^2)_i} = \lambda(\sqrt{A^2})_i = \lambda(A)_i > 0, & 1 \leq i \leq \rho \\ 0, & \rho < i \leq \eta \end{cases} \quad (1729)$$

A positive semidefinite matrix, having diagonalization  $A = S\Lambda S^T$  (1711) and full singular value decomposition  $A = U\Sigma Q^T$  (1726), simply relates the two

$$A = S\Lambda S^T = U\Sigma Q^T \quad (1730)$$

by direct correspondence; *id est*,  $S = U = Q$ ,  $\Lambda = \Sigma$ .

### A.6.3.2 SVD of symmetric matrices

From (1719) and (1715) for  $A = A^T$ , more generally,

$$\sigma(A)_i = \begin{cases} \sqrt{\lambda(A^2)_i} = \lambda(\sqrt{A^2})_i = |\lambda(A)_i| > 0, & 1 \leq i \leq \rho \\ 0, & \rho < i \leq \eta \end{cases} \quad (1731)$$

For any  $A \in \mathbb{R}^{m \times n}$ ,  $A^T A$  is (symmetric) positive semidefinite:

$$\sigma(A^T A)_i = \begin{cases} \lambda(A^T A)_i > 0, & 1 \leq i \leq \rho \\ 0, & \rho < i \leq n \end{cases} \quad (1732)$$

#### A.6.3.2.1 Definition. Step function. (confer §4.3.2.0.1)

Define the signum-like quasilinear function  $\psi : \mathbb{R}^n \rightarrow \mathbb{R}^n$  that takes value 1 corresponding to a 0-valued entry in its argument:

$$\psi(a) \triangleq \left[ \lim_{x_i \rightarrow a_i} \frac{x_i}{|x_i|} \right] = \begin{cases} 1, & a_i \geq 0 \\ -1, & a_i < 0 \end{cases}, \quad i = 1 \dots n \in \mathbb{R}^n \quad (1733)$$

Unlike  $\text{sgn}()$ ,  $\psi$  is not an odd function;  $\psi(-a) \neq -\psi(a)$  because of 0 handling.  $\triangle$

Eigenvalue signs of a symmetric matrix having diagonalization  $A = S\Lambda S^T$  (1711) can be absorbed either into real  $U$  or real  $Q$  from the full SVD; [391, p.34] (confer §C.4.2.1)

$$A = S\Lambda S^T = S\delta(\psi(\delta(\Lambda)))|\Lambda|S^T \triangleq U\Sigma Q^T \in \mathbb{S}^n \quad (1734)$$

or

$$A = S\Lambda S^T = S|\Lambda|\delta(\psi(\delta(\Lambda)))S^T \triangleq U\Sigma Q^T \in \mathbb{S}^n \quad (1735)$$

where matrix of singular values  $\Sigma = |\Lambda|$  denotes entrywise absolute value of diagonal eigenvalue matrix  $\Lambda$ .

#### A.6.4 Pseudoinverse by SVD

Matrix pseudoinverse (§E) is nearly synonymous with singular value decomposition because of the elegant expression, given  $A = U\Sigma Q^T \in \mathbb{R}^{m \times n}$

$$A^\dagger = Q\Sigma^{\dagger T}U^T \in \mathbb{R}^{n \times m} \quad (1736)$$

that applies to all three flavors of SVD, where  $\Sigma^\dagger$  simply inverts nonzero entries of matrix  $\Sigma$ .

Given symmetric matrix  $A \in \mathbb{S}^n$  and its diagonalization  $A = S\Lambda S^T$  (§A.5.1), its pseudoinverse simply inverts all nonzero eigenvalues:

$$A^\dagger = S\Lambda^\dagger S^T \in \mathbb{S}^n \quad (1737)$$

### A.7 Zeros

#### A.7.1 norm zero

For any given norm, by definition,

$$\|x\|_\ell = 0 \Leftrightarrow x = \mathbf{0} \quad (1738)$$

Consequently, a generally nonconvex constraint in  $x$  like  $\|Ax - b\| = \kappa$  becomes convex when  $\kappa = 0$ .

#### A.7.2 0 entry

If a positive semidefinite matrix  $A = [A_{ij}] \in \mathbb{R}^{n \times n}$  has a 0 entry  $A_{ii}$  on its main diagonal, then  $A_{ij} + A_{ji} = 0 \forall j$ . [303, §1.3.1]

Any symmetric positive semidefinite matrix  $A \in \mathbb{S}^n$ , having a 0 entry on its main diagonal, must be  $\mathbf{0}$  along the entire row and column to which that 0 entry belongs. [181, §4.2.8] [228, §7.1 prob.2] From which it follows:

$$\delta(A) = \mathbf{0} \Leftrightarrow A = \mathbf{0} \quad (1739)$$

$$\text{tr}(A) = 0 \Leftrightarrow A = \mathbf{0} \quad (1740)$$

##### A.7.2.0.1 Exercise. Positive semidefinite matrix diagonal zero.

Which Schur complement condition demands multiple 0 entries in submatrix  $B$  when there is a single 0 entry on the main diagonal of submatrix  $A$  in partitioned positive semidefinite matrix  $G$  in (1663)? Having made that determination, can one show consequent necessity for a zero row and column in  $G$  simply by repartitioning?

In the same regard, which condition principally governs case  $C = \mathbf{0}$ ? Prove that  $B(I - CC^\dagger) = \mathbf{0}$  is a necessary condition. ▼

#### A.7.3 0 eigenvalues theorem

This theorem is simple, powerful, and widely applicable:

##### A.7.3.0.1 Theorem. Number of 0 eigenvalues.

- For any matrix  $A \in \mathbb{R}^{m \times n}$

$$\text{rank}(A) + \dim \mathcal{N}(A) = n \quad (1741)$$

by conservation of dimension. [228, §0.4.4]

- For any square matrix  $A \in \mathbb{R}^{m \times m}$ , number of 0 eigenvalues is at least equal to  $\dim \mathcal{N}(A)$ ;

$$\dim \mathcal{N}(A) \leq \text{number of 0 eigenvalues} \leq m \quad (1742)$$

All eigenvectors, corresponding to those 0 eigenvalues, belong to  $\mathcal{N}(A)$ . [\[368, §5.1\]](#)

- For diagonalizable matrix  $A$  ([§A.5](#)), the number of 0 eigenvalues is precisely  $\dim \mathcal{N}(A)$  while the corresponding eigenvectors span  $\mathcal{N}(A)$ . Real and imaginary parts of the eigenvectors remaining span  $\mathcal{R}(A)$ .

(TRANSPOSE.)

- Likewise, for any matrix  $A \in \mathbb{R}^{m \times n}$

$$\text{rank}(A^T) + \dim \mathcal{N}(A^T) = m \quad (1743)$$

- For any square  $A \in \mathbb{R}^{m \times m}$ , number of 0 eigenvalues is at least equal to  $\dim \mathcal{N}(A^T) = \dim \mathcal{N}(A)$ . All left-eigenvectors (eigenvectors of  $A^T$ ), corresponding to those 0 eigenvalues, belong to  $\mathcal{N}(A^T)$ .
- For diagonalizable  $A$ , number of 0 eigenvalues is precisely  $\dim \mathcal{N}(A^T)$  while the corresponding left-eigenvectors span  $\mathcal{N}(A^T)$ . Real and imaginary parts of the left-eigenvectors remaining span  $\mathcal{R}(A^T)$ .  $\diamond$

**Proof.** First we show, for a diagonalizable matrix, the number of 0 eigenvalues is precisely the dimension of its nullspace while the eigenvectors corresponding to those 0 eigenvalues span the nullspace:

Any diagonalizable matrix  $A \in \mathbb{R}^{m \times m}$  must possess a complete set of linearly independent eigenvectors. If  $A$  is full-rank (invertible), then all  $m = \text{rank}(A)$  eigenvalues are nonzero. [\[368, §5.1\]](#)

Suppose  $\text{rank}(A) < m$ . Then  $\dim \mathcal{N}(A) = m - \text{rank}(A)$ . Thus there is a set of  $m - \text{rank}(A)$  linearly independent vectors spanning  $\mathcal{N}(A)$ . Each of those can be an eigenvector associated with a 0 eigenvalue because  $A$  is diagonalizable  $\Leftrightarrow \exists m$  linearly independent eigenvectors. [\[368, §5.2\]](#) Eigenvectors of a real matrix corresponding to 0 eigenvalues must be real. [A.19](#) Thus  $A$  has at least  $m - \text{rank}(A)$  eigenvalues equal to 0.

Now suppose  $A$  has more than  $m - \text{rank}(A)$  eigenvalues equal to 0. Then there are more than  $m - \text{rank}(A)$  linearly independent eigenvectors associated with 0 eigenvalues, and each of those eigenvectors must be in  $\mathcal{N}(A)$ . Thus there are more than  $m - \text{rank}(A)$  linearly independent vectors in  $\mathcal{N}(A)$ ; a contradiction.

Diagonalizable  $A$  therefore has  $\text{rank}(A)$  nonzero eigenvalues and exactly  $m - \text{rank}(A)$  eigenvalues equal to 0 whose corresponding eigenvectors span  $\mathcal{N}(A)$ .

By similar argument, the left-eigenvectors corresponding to 0 eigenvalues span  $\mathcal{N}(A^T)$ .

---

[A.18](#) We take as given the well-known fact that the number of 0 eigenvalues cannot be less than dimension of the nullspace. We offer an example of the converse:

$$A = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

$\dim \mathcal{N}(A) = 2$ ,  $\lambda(A) = [0 \ 0 \ 0 \ 1]^T$ ; three eigenvectors in the nullspace but only two are independent. The right-hand side of (1742) is tight for nonzero matrices; e.g., ([§B.1](#)) dyad  $uv^T \in \mathbb{R}^{m \times m}$  has  $m$  0-eigenvalues when  $u \in v^\perp$ .

[A.19 Proof.](#) Let  $*$  denote complex conjugation. Suppose  $A = A^*$  and  $As_i = \mathbf{0}$ . Then  $s_i = s_i^* \Rightarrow As_i = As_i^* \Rightarrow As_i^* = \mathbf{0}$ . Conversely,  $As_i^* = \mathbf{0} \Rightarrow As_i = As_i^* \Rightarrow s_i = s_i^*$ .  $\diamond$

Next we show when  $A$  is diagonalizable, the real and imaginary parts of its eigenvectors (corresponding to nonzero eigenvalues) span  $\mathcal{R}(A)$ :

The (right-)eigenvectors of a diagonalizable matrix  $A \in \mathbb{R}^{m \times m}$  are linearly independent if and only if the left-eigenvectors are. So, matrix  $A$  has a representation in terms of its right- and left-eigenvectors; from the diagonalization (1699), assuming 0 eigenvalues are ordered last,

$$A = \sum_{\substack{i=1 \\ \lambda_i \neq 0}}^m \lambda_i s_i w_i^T = \sum_{\substack{i=1 \\ \lambda_i \neq 0}}^{k \leq m} \lambda_i s_i w_i^T \quad (1744)$$

From the *linearly independent dyads theorem* (§B.1.1.0.2), the dyads  $\{s_i w_i^T\}$  must be independent because each set of eigenvectors are; hence  $\text{rank } A = k$ , the number of nonzero eigenvalues. Complex eigenvectors and eigenvalues are common for real matrices, and must come in complex conjugate pairs for the summation to remain real. Assume that conjugate pairs of eigenvalues appear in sequence. Given any particular conjugate pair from (1744), we get the partial summation

$$\begin{aligned} \lambda_i s_i w_i^T + \lambda_i^* s_i^* w_i^{*T} &= 2 \operatorname{re}(\lambda_i s_i w_i^T) \\ &= 2(\operatorname{re} s_i \operatorname{re}(\lambda_i w_i^T) - \operatorname{im} s_i \operatorname{im}(\lambda_i w_i^T)) \end{aligned} \quad (1745)$$

where<sup>A.20</sup>  $\lambda_i^* \triangleq \lambda_{i+1}$ ,  $s_i^* \triangleq s_{i+1}$ , and  $w_i^* \triangleq w_{i+1}$ . Then (1744) is equivalently written

$$A = 2 \sum_{\substack{i \\ \lambda \in \mathbb{C} \\ \lambda_i \neq 0}} \operatorname{re} s_{2i} \operatorname{re}(\lambda_{2i} w_{2i}^T) - \operatorname{im} s_{2i} \operatorname{im}(\lambda_{2i} w_{2i}^T) + \sum_{\substack{j \\ \lambda \in \mathbb{R} \\ \lambda_j \neq 0}} \lambda_j s_j w_j^T \quad (1746)$$

The summation (1746) shows:  $A$  is a linear combination of real and imaginary parts of its (right-)eigenvectors corresponding to nonzero eigenvalues. The  $k$  vectors  $\{\operatorname{re} s_i \in \mathbb{R}^m, \operatorname{im} s_i \in \mathbb{R}^m \mid \lambda_i \neq 0, i \in \{1 \dots m\}\}$  must therefore span the range of diagonalizable matrix  $A$ .

The argument is similar regarding span of the left-eigenvectors. ♦

#### A.7.4 0 trace and matrix product

For  $X, A \in \mathbb{R}_+^{M \times N}$  (39)

$$\operatorname{tr}(X^T A) = 0 \Leftrightarrow X \circ A = A \circ X = \mathbf{0} \quad (1747)$$

For  $X, A \in \mathbb{S}_+^M$  [35, §2.6.1 exer.2.8] [400, §3.1]

$$\operatorname{tr}(XA) = 0 \Leftrightarrow XA = AX = \mathbf{0} \quad (1748)$$

**Proof.** ( $\Leftarrow$ ) Suppose  $XA = AX = \mathbf{0}$ . Then  $\operatorname{tr}(XA) = 0$  is obvious.

( $\Rightarrow$ ) Suppose  $\operatorname{tr}(XA) = 0$ .  $\operatorname{tr}(XA) = \operatorname{tr}(\sqrt{A} X \sqrt{A})$  whose argument is positive semidefinite by Corollary A.3.1.0.5. Trace of any square matrix is equivalent to the sum of its eigenvalues. Eigenvalues of a positive semidefinite matrix can total 0 if and only if each and every nonnegative eigenvalue is 0. The only positive semidefinite matrix, having all 0 eigenvalues, resides at the origin; (*confer*(1772)) *id est*,

$$\sqrt{A} X \sqrt{A} = (\sqrt{X} \sqrt{A})^T \sqrt{X} \sqrt{A} = \mathbf{0} \quad (1749)$$

---

<sup>A.20</sup>Complex conjugate of  $w$  is denoted  $w^*$ . Conjugate transpose is denoted  $w^H = w^{*T}$ .

implying  $\sqrt{X}\sqrt{A} = \mathbf{0}$  which in turn implies  $\sqrt{X}(\sqrt{X}\sqrt{A})\sqrt{A} = XA = \mathbf{0}$ . Arguing similarly yields  $AX = \mathbf{0}$ .  $\blacklozenge$

Diagonalizable matrices  $A$  and  $X$  are *simultaneously diagonalizable* if and only if they are commutative under multiplication; [228, §1.3.12] *id est*, iff they share a complete set of eigenvectors.

#### A.7.4.0.1 Example. An equivalence in nonisomorphic spaces.

Identity (1748) leads to an unusual equivalence relating convex geometry to traditional linear algebra: The convex sets, given  $A \succeq 0$

$$\{X \mid \langle X, A \rangle = 0\} \cap \{X \succeq 0\} \equiv \{X \mid \mathcal{N}(X) \supseteq \mathcal{R}(A)\} \cap \{X \succeq 0\} \quad (1750)$$

(one expressed in terms of a hyperplane, the other in terms of nullspace and range) are equivalent only when symmetric matrix  $A$  is positive semidefinite.

We might apply this equivalence to the geometric center subspace, for example,

$$\begin{aligned} \mathbb{S}_c^M &= \{Y \in \mathbb{S}^M \mid Y\mathbf{1} = \mathbf{0}\} \\ &= \{Y \in \mathbb{S}^M \mid \mathcal{N}(Y) \supseteq \mathbf{1}\} = \{Y \in \mathbb{S}^M \mid \mathcal{R}(Y) \subseteq \mathcal{N}(\mathbf{1}^T)\} \end{aligned} \quad (2196)$$

from which we derive (*confer* (1139))

$$\mathbb{S}_c^M \cap \mathbb{S}_+^M \equiv \{X \succeq 0 \mid \langle X, \mathbf{1}\mathbf{1}^T \rangle = 0\} \quad (1751)$$

$\square$

#### A.7.5 Zero definite

The domain over which an arbitrary real matrix  $A \in \mathbb{R}^{M \times M}$  is zero definite can exceed its left and right nullspaces; *e.g.*, [455, §3.2 prob.5]

$$\{x \mid x^T A x = 0\} = \mathbb{R}^M \Leftrightarrow A^T = -A \quad (1752)$$

whereas

$$\{x \mid x^H A x = 0\} = \mathbb{C}^M \Leftrightarrow A = \mathbf{0} \quad (1753)$$

For any positive semidefinite matrix  $A \in \mathbb{R}^{M \times M}$  (for  $A + A^T \succeq 0$ )

$$\{x \mid x^T A x = 0\} = \mathcal{N}(A + A^T) \quad (1754)$$

because  $\exists R \ni A + A^T = R^T R$ ,  $\|Rx\| = 0 \Leftrightarrow Rx = \mathbf{0}$ , and  $\mathcal{N}(A + A^T) = \mathcal{N}(R)$ . Then given any particular vector  $x_p$ ,  $x_p^T A x_p = 0 \Leftrightarrow x_p \in \mathcal{N}(A + A^T)$ .

For any positive definite matrix  $A \in \mathbb{R}^{M \times M}$  (for  $A + A^T \succ 0$ )

$$\{x \mid x^T A x = 0\} = \mathbf{0} \quad (1755)$$

#### A.7.5.0.1 Example. Zero definiteness.

The positive semidefinite matrix

$$A = \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \quad (1756)$$

has no nullspace. Yet

$$\{x \mid x^T A x = 0\} = \{x \mid \mathbf{1}^T x = 0\} \subset \mathbb{R}^2 \quad (1757)$$

which is the nullspace of the symmetrized matrix. Symmetric matrices are not spared from the excess; *videlicet*,

$$B = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \quad (1758)$$

has eigenvalues  $\{-1, 3\}$ , no nullspace, but is zero definite on [A.21](#)

$$\mathcal{X} \triangleq \{x \in \mathbb{R}^2 \mid x_2 = (-2 \pm \sqrt{3})x_1\} \quad (1759)$$

□

#### A.7.5.0.2 Proposition. (Sturm/Zhang) Dyad-decompositions. [374, §5.2]

Let positive semidefinite matrix  $X \in \mathbb{S}_+^M$  have rank  $\rho$ . Then, given symmetric matrix  $A \in \mathbb{S}^M$ ,  $\langle A, X \rangle = 0$  if and only if there exists a dyad-decomposition

$$X = \sum_{j=1}^{\rho} x_j x_j^T \quad (1760)$$

satisfying

$$\langle A, x_j x_j^T \rangle = 0 \text{ for each and every } j \in \{1 \dots \rho\} \quad (1761)$$

◊

The dyad-decomposition of  $X$  proposed is generally not that obtained from a standard diagonalization by eigenvalue decomposition, unless  $\rho=1$  or the given matrix  $A$  is simultaneously diagonalizable ([§A.7.4](#)) with  $X$ . That means, elemental dyads in decomposition (1760) constitute a generally nonorthogonal set. Sturm & Zhang give a simple procedure for constructing the dyad-decomposition [430] where matrix  $A$  may be regarded as a parameter.

#### A.7.5.0.3 Example. Dyad.

The dyad  $uv^T \in \mathbb{R}^{M \times M}$  ([§B.1](#)) is zero definite on all  $x$  for which either  $x^T u = 0$  or  $x^T v = 0$ ;

$$\{x \mid x^T uv^T x = 0\} = \{x \mid x^T u = 0\} \cup \{x \mid v^T x = 0\} \quad (1762)$$

*id est*, on  $u^\perp \cup v^\perp$ . Symmetrizing the dyad does not change the outcome:

$$\{x \mid x^T (uv^T + vu^T)x/2 = 0\} = \{x \mid x^T u = 0\} \cup \{x \mid v^T x = 0\} \quad (1763)$$

□

---

[A.21](#) These two lines represent the limit in the union of two generally distinct hyperbolae; *id est*, for matrix  $B$  and set  $\mathcal{X}$  as defined

$$\lim_{\varepsilon \rightarrow 0^+} \{x \in \mathbb{R}^2 \mid x^T B x = \varepsilon\} = \mathcal{X}$$

# Appendix B

## Simple matrices

*Mathematicians also attempted to develop algebra of vectors but there was no natural definition of the product of two vectors that held in arbitrary dimensions. The first vector algebra that involved a noncommutative vector product (that is,  $v \times w$  need not equal  $w \times v$ ) was proposed by Hermann Grassmann in his book Ausdehnungslehre (1844). Grassmann's text also introduced the product of a column matrix and a row matrix, which resulted in what is now called a simple or a rank-one matrix. In the late 19th century the American mathematical physicist Willard Gibbs published his famous treatise on vector analysis. In that treatise Gibbs represented general matrices, which he called dyadics, as sums of simple matrices, which Gibbs called dyads. Later the physicist P. A. M. Dirac introduced the term "bra-ket" for what we now call the scalar product of a "bra" (row) vector times a "ket" (column) vector and the term "ket-bra" for the product of a ket times a bra, resulting in what we now call a simple matrix, as above. Our convention of identifying column matrices and vectors was introduced by physicists in the 20th century.*

—Suddhendu Biswas [51, p.2]

### B.1 Rank-one matrix (dyad)

Any matrix formed from the unsigned outer product of two vectors,

$$\Psi = uv^T \in \mathbb{R}^{M \times N} \quad (1764)$$

where  $u \in \mathbb{R}^M$  and  $v \in \mathbb{R}^N$ , is rank-one and called a *dyad*. Conversely, any rank-one matrix must have the form  $\Psi$ . [228, prob.1.4.1] Product  $-uv^T$  is a *negative dyad*. For matrix products  $AB^T$ , in general, we have

$$\mathcal{R}(AB^T) \subseteq \mathcal{R}(A), \quad \mathcal{N}(AB^T) \supseteq \mathcal{N}(B^T) \quad (1765)$$

with equality when  $B = A$  [368, §3.3, §3.6]<sup>B.1</sup> or respectively when  $B$  is invertible and  $\mathcal{N}(A) = \mathbf{0}$ . Yet for all nonzero dyads we have

$$\mathcal{R}(uv^T) = \mathcal{R}(u), \quad \mathcal{N}(uv^T) = \mathcal{N}(v^T) \equiv v^\perp \quad (1766)$$

---

**B.1 Proof.**  $\mathcal{R}(AA^T) \subseteq \mathcal{R}(A)$  is obvious.

$$\begin{aligned} \mathcal{R}(AA^T) &= \{AA^Ty \mid y \in \mathbb{R}^m\} \\ &\supseteq \{AA^Ty \mid A^Ty \in \mathcal{R}(A^T)\} = \mathcal{R}(A) \text{ by (144)} \end{aligned}$$

♦

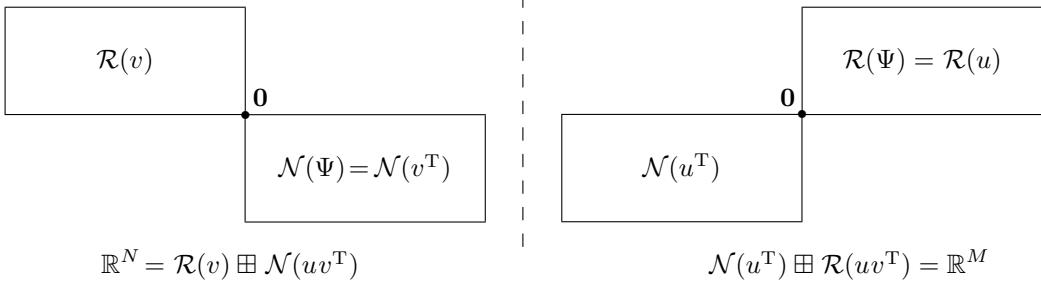


Figure 181: The four fundamental subspaces [370, §3.6] for any dyad  $\Psi = uv^T \in \mathbb{R}^{M \times N}$ :  $\mathcal{R}(v) \perp \mathcal{N}(\Psi)$  &  $\mathcal{N}(u^T) \perp \mathcal{R}(\Psi)$ .  $\Psi(x) \triangleq uv^T x$  is a linear mapping from  $\mathbb{R}^N$  to  $\mathbb{R}^M$ . Map from  $\mathcal{R}(v)$  to  $\mathcal{R}(u)$  is bijective. [368, §3.1]

where  $\dim v^\perp = N - 1$ .

It is obvious that a dyad can be  $\mathbf{0}$  only when  $u$  or  $v$  is  $\mathbf{0}$ ;

$$\Psi = uv^T = \mathbf{0} \Leftrightarrow u = \mathbf{0} \text{ or } v = \mathbf{0} \quad (1767)$$

The matrix 2-norm for  $\Psi$  is equivalent to Frobenius' norm;

$$\|\Psi\|_2 = \sigma_1 = \|uv^T\|_F = \|uv^T\|_2 = \|u\| \|v\| \quad (1768)$$

When  $u$  and  $v$  are normalized, the pseudoinverse is the transposed dyad. Otherwise,

$$\Psi^\dagger = (uv^T)^\dagger = \frac{vu^T}{\|u\|^2 \|v\|^2} \quad (1769)$$

When dyad  $uv^T \in \mathbb{R}^{N \times N}$  is square,  $uv^T$  has at least  $N - 1$  0-eigenvalues and corresponding eigenvectors spanning  $v^\perp$ . The remaining eigenvector  $u$  spans the range of  $uv^T$  with corresponding eigenvalue

$$\lambda = v^T u = \text{tr}(uv^T) \in \mathbb{R} \quad (1770)$$

Determinant is a product of the eigenvalues; so, it is always true that

$$\det \Psi = \det(uv^T) = 0 \quad (1771)$$

When  $\lambda = 1$ , the square dyad is a nonorthogonal projector projecting on its range ( $\Psi^2 = \Psi$ , §E.6); a *projector dyad*. It is quite possible that  $u \in v^\perp$  making the remaining eigenvalue instead 0; <sup>B.2</sup>  $\lambda = 0$  together with the first  $N - 1$  0-eigenvalues; *id est*, it is possible  $uv^T$  were nonzero while all its eigenvalues are 0. The matrix

$$\begin{bmatrix} 1 \\ -1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix} \quad (1772)$$

for example, has two 0-eigenvalues. In other words, eigenvector  $u$  may simultaneously be a member of the nullspace and range of the dyad. The explanation is, simply, because  $u$  and  $v$  share the same dimension,  $\dim u = M = \dim v = N$ :

---

<sup>B.2</sup>A dyad is not always diagonalizable (§A.5) because its eigenvectors are not necessarily independent.

**Proof.** Figure 181 sees the four fundamental subspaces for the dyad. Linear operator  $\Psi : \mathbb{R}^N \rightarrow \mathbb{R}^M$  provides a map between vector spaces that remain distinct when  $M = N$ ;

$$\begin{aligned} u &\in \mathcal{R}(uv^T) \\ u &\in \mathcal{N}(uv^T) \Leftrightarrow v^Tu = 0 \\ \mathcal{R}(uv^T) \cap \mathcal{N}(uv^T) &= \emptyset \end{aligned} \quad (1773)$$

◆

### B.1.0.1 rank-one modification

For  $A \in \mathbb{R}^{M \times N}$ ,  $x \in \mathbb{R}^N$ ,  $y \in \mathbb{R}^M$ , and  $y^TAx \neq 0$  [233, §2.1]<sup>B.3</sup>

$$\text{rank}\left(A - \frac{Axy^TA}{y^TAx}\right) = \text{rank}(A) - 1 \quad (1775)$$

Given nonsingular matrix  $A \in \mathbb{R}^{N \times N}$   $\exists 1 + v^TA^{-1}u \neq 0$ , [174, §4.11.2] [248, App.6] [455, §2.3 prob.16] (Sherman-Morrison-Woodbury)<sup>B.4</sup>

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u} \quad (1776)$$

### B.1.0.2 dyad symmetry

In the specific circumstance that  $v = u$ , then  $uu^T \in \mathbb{R}^{N \times N}$  is symmetric, rank-one, and positive semidefinite having exactly  $N - 1$  0-eigenvalues. In fact, (Theorem A.3.1.0.7)

$$uv^T \succeq 0 \Leftrightarrow v = u \quad (1777)$$

and the remaining eigenvalue is almost always positive;

$$\lambda = u^Tu = \text{tr}(uu^T) > 0 \text{ unless } u = 0 \quad (1778)$$

When  $\lambda = 1$ , the dyad becomes an orthogonal projector.

Matrix

$$\begin{bmatrix} \Psi & u \\ u^T & 1 \end{bmatrix} \quad (1779)$$

for example, is rank-1 positive semidefinite if and only if  $\Psi = uu^T$ .

## B.1.1 Dyad independence

Now we consider a sum of dyads like (1764) as encountered in diagonalization and singular value decomposition:

$$\mathcal{R}\left(\sum_{i=1}^k s_i w_i^T\right) = \sum_{i=1}^k \mathcal{R}(s_i w_i^T) = \sum_{i=1}^k \mathcal{R}(s_i) \Leftarrow w_i \forall i \text{ are l.i.} \quad (1780)$$

<sup>B.3</sup>This rank-one modification formula has a Schur progenitor, in the symmetric case:

$$\begin{array}{ll} \underset{c}{\text{minimize}} & c \\ \text{subject to} & \begin{bmatrix} A & Ax \\ y^TA & c \end{bmatrix} \succeq 0 \end{array} \quad (1774)$$

has analytical solution by (1663b):  $c \geq y^TAx$ . Difference  $A - \frac{Axy^TA}{y^TAx}$  comes from (1663c). Rank modification is provable via Theorem A.4.0.1.3.

<sup>B.4</sup> which is unstable numerically, by **Saunders' reckoning**: “We know how to update matrix factorizations reliably (even if the input or output matrix is singular) but, in general, there’s no stable way to update a matrix inverse.”

range of summation is the vector sum of ranges.<sup>B.5</sup> (Theorem B.1.1.1) Under the assumption the dyads are linearly independent (l.i.), then vector sums are unique (p.624): for  $\{w_i\}$  l.i. and  $\{s_i\}$  l.i.

$$\mathcal{R}\left(\sum_{i=1}^k s_i w_i^T\right) = \mathcal{R}(s_1 w_1^T) \oplus \dots \oplus \mathcal{R}(s_k w_k^T) = \mathcal{R}(s_1) \oplus \dots \oplus \mathcal{R}(s_k) \quad (1781)$$

**B.1.1.0.1 Definition.** *Linearly independent dyads.* [238, p.29 thm.11] [376, p.2]  
The set of  $k$  dyads

$$\{s_i w_i^T \mid i=1 \dots k\} \quad (1782)$$

where  $s_i \in \mathbb{C}^M$  and  $w_i \in \mathbb{C}^N$ , is said to be linearly independent iff

$$\text{rank}\left(SW^T \triangleq \sum_{i=1}^k s_i w_i^T\right) = k \quad (1783)$$

where  $S \triangleq [s_1 \dots s_k] \in \mathbb{C}^{M \times k}$  and  $W \triangleq [w_1 \dots w_k] \in \mathbb{C}^{N \times k}$ .  $\triangle$

Dyad independence does not preclude existence of a nullspace  $\mathcal{N}(SW^T)$ , as defined, nor does it imply  $SW^T$  were full-rank. In absence of assumption of independence, generally,  $\text{rank } SW^T \leq k$ . Conversely, any rank- $k$  matrix can be written in the form  $SW^T$  by singular value decomposition. (§A.6)

**B.1.1.0.2 Theorem.** *Linearly independent (l.i.) dyads.*

Vectors  $\{s_i \in \mathbb{C}^M, i=1 \dots k\}$  are l.i. and vectors  $\{w_i \in \mathbb{C}^N, i=1 \dots k\}$  are l.i. if and only if dyads  $\{s_i w_i^T \in \mathbb{C}^{M \times N}, i=1 \dots k\}$  are l.i.  $\diamond$

**Proof.** Linear independence of  $k$  dyads is identical to definition (1783).

( $\Rightarrow$ ) Suppose  $\{s_i\}$  and  $\{w_i\}$  are each linearly independent sets. Invoking Sylvester's rank inequality, [228, §0.4] [455, §2.4]

$$\text{rank } S + \text{rank } W - k \leq \text{rank}(SW^T) \leq \min\{\text{rank } S, \text{rank } W\} (\leq k) \quad (1784)$$

Then  $k \leq \text{rank}(SW^T) \leq k$  which implies the dyads are independent.

( $\Leftarrow$ ) Conversely, suppose  $\text{rank}(SW^T) = k$ . Then

$$k \leq \min\{\text{rank } S, \text{rank } W\} \leq k \quad (1785)$$

implying the vector sets are each independent.  $\spadesuit$

### B.1.1.1 Biorthogonality condition, Range and Nullspace of Sum

Dyads characterized by biorthogonality condition  $W^T S = I$  are independent; *id est*, for  $S \in \mathbb{C}^{M \times k}$  and  $W \in \mathbb{C}^{N \times k}$ , if  $W^T S = I$  then  $\text{rank}(SW^T) = k$  by the *linearly independent dyads theorem* because (confer §E.1.1)

$$W^T S = I \Rightarrow \text{rank } S = \text{rank } W = k \leq M = N \quad (1786)$$

To see that, we need only show:  $\mathcal{N}(S) = \mathbf{0} \Leftrightarrow \exists B \ni BS = I$ .<sup>B.6</sup>

( $\Leftarrow$ ) Assume  $BS = I$ . Then  $\mathcal{N}(BS) = \mathbf{0} = \{x \mid BSx = \mathbf{0}\} \supseteq \mathcal{N}(S)$ . (1765)

<sup>B.5</sup>Move of range  $\mathcal{R}$  to inside summation admitted by linear independence of  $\{w_i\}$ .

<sup>B.6</sup>Left inverse is not unique, in general.

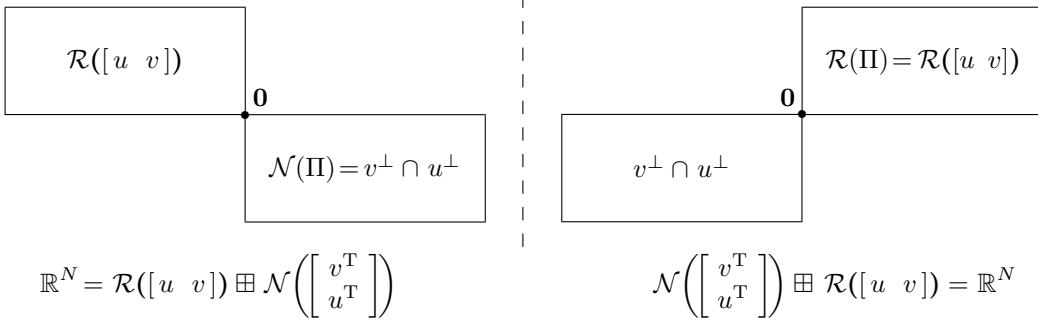


Figure 182: The four fundamental subspaces [370, §3.6] for doublet  $\Pi = uv^T + vu^T \in \mathbb{S}^N$ .  $\Pi(x) = (uv^T + vu^T)x$  is a linear bijective mapping from  $\mathcal{R}([u v])$  to  $\mathcal{R}([u v])$ .

$(\Rightarrow)$  If  $\mathcal{N}(S) = \mathbf{0}$  then  $S$  must be full-rank thin-or-square.

$$\therefore \exists A, B, C \ni \begin{bmatrix} B \\ C \end{bmatrix} [S A] = I \text{ (id est, } [S A] \text{ is invertible)} \Rightarrow BS = I.$$

Left inverse  $B$  is given as  $W^T$  here. Because of reciprocity with  $S$ , it immediately follows:  $\mathcal{N}(W) = \mathbf{0} \Leftrightarrow \exists S \ni S^T W = I$ .  $\blacklozenge$

Dyads produced by diagonalization, for example, are independent because of their inherent biorthogonality. (§A.5.0.3) The converse is generally false; *id est*, linearly independent dyads are not necessarily biorthogonal.

#### B.1.1.1.1 Theorem. Nullspace and range of dyad sum.

Given a sum of dyads represented by  $SW^T$  where  $S \in \mathbb{C}^{M \times k}$  and  $W \in \mathbb{C}^{N \times k}$

$$\begin{aligned} \mathcal{N}(SW^T) &= \mathcal{N}(W^T) \Leftarrow \exists B \ni BS = I \\ \mathcal{R}(SW^T) &= \mathcal{R}(S) \Leftarrow \exists Z \ni W^T Z = I \end{aligned} \quad (1787)$$

$\diamond$

**Proof.**  $(\Rightarrow)$   $\mathcal{N}(SW^T) \supseteq \mathcal{N}(W^T)$  and  $\mathcal{R}(SW^T) \subseteq \mathcal{R}(S)$  are obvious.

$(\Leftarrow)$  Assume existence of a left inverse  $B \in \mathbb{R}^{k \times N}$  and a right inverse  $Z \in \mathbb{R}^{N \times k}$ . B.7

$$\mathcal{N}(SW^T) = \{x \mid SW^T x = \mathbf{0}\} \subseteq \{x \mid BSW^T x = \mathbf{0}\} = \mathcal{N}(W^T) \quad (1788)$$

$$\mathcal{R}(SW^T) = \{SW^T x \mid x \in \mathbb{R}^N\} \supseteq \{SW^T Z y \mid Z y \in \mathbb{R}^N\} = \mathcal{R}(S) \quad (1789)$$

$\blacklozenge$

## B.2 Doublet

Consider a sum of two linearly independent square dyads, one a transposition of the other:

$$\Pi = uv^T + vu^T = \begin{bmatrix} u & v \end{bmatrix} \begin{bmatrix} v^T \\ u^T \end{bmatrix} = SW^T \in \mathbb{S}^N \quad (1790)$$

where  $u, v \in \mathbb{R}^N$ . Like the dyad, a doublet can be  $\mathbf{0}$  only when  $u$  or  $v$  is  $\mathbf{0}$ ;

$$\Pi = uv^T + vu^T = \mathbf{0} \Leftrightarrow u = \mathbf{0} \text{ or } v = \mathbf{0} \quad (1791)$$

By assumption of independence, a nonzero doublet has two nonzero eigenvalues

$$\lambda_1 \triangleq u^T v + \|uv^T\|, \quad \lambda_2 \triangleq u^T v - \|uv^T\| \quad (1792)$$

B.7 By counterexample, the theorem's converse cannot be true; e.g.,  $S = W = [\mathbf{1} \ \mathbf{0}]$ .

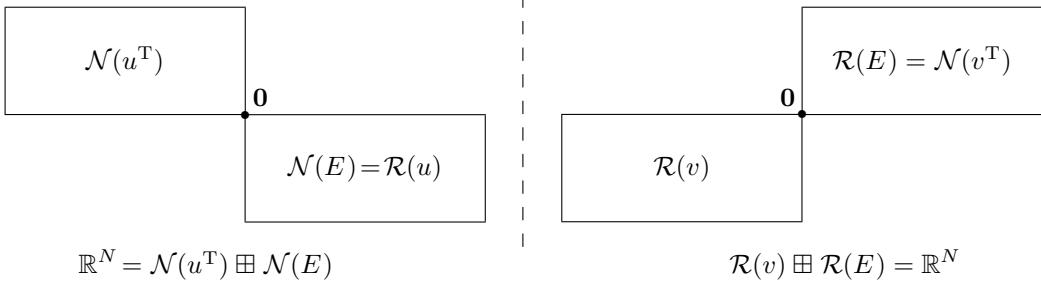


Figure 183:  $v^T u = 1/\zeta$ . The four fundamental subspaces [370, §3.6] for elementary matrix  $E$  as a linear mapping  $E(x) = \left( I - \frac{uv^T}{v^T u} \right) x$ .

where  $\lambda_1 > 0 > \lambda_2$ , with corresponding eigenvectors

$$x_1 \triangleq \frac{u}{\|u\|} + \frac{v}{\|v\|}, \quad x_2 \triangleq \frac{u}{\|u\|} - \frac{v}{\|v\|} \quad (1793)$$

spanning the doublet range. Eigenvalue  $\lambda_1$  cannot be 0 unless  $u$  and  $v$  have opposing directions, but that is antithetical since then the dyads would no longer be independent. Eigenvalue  $\lambda_2$  is 0 if and only if  $u$  and  $v$  share the same direction, again antithetical. Generally we have  $\lambda_1 > 0$  and  $\lambda_2 < 0$ , so  $\Pi$  is indefinite.

By the *nullspace and range of dyad sum theorem*, doublet  $\Pi$  has  $N-2$  zero-eigenvalues remaining and corresponding eigenvectors spanning  $\mathcal{N}\left(\begin{bmatrix} v^T \\ u^T \end{bmatrix}\right)$ . We therefore have

$$\mathcal{R}(\Pi) = \mathcal{R}([u \ v]), \quad \mathcal{N}(\Pi) = v^\perp \cap u^\perp \quad (1794)$$

of respective dimension 2 and  $N-2$ . (Figure 182)

### B.3 Elementary matrix

A matrix of the form

$$E = I - \zeta uv^T \in \mathbb{R}^{N \times N} \quad (1795)$$

where  $\zeta \in \mathbb{R}$  is finite and  $u, v \in \mathbb{R}^N$ , is called an *elementary matrix* or *rank-one modification of the Identity*. [230] Any elementary matrix in  $\mathbb{R}^{N \times N}$  has  $N-1$  eigenvalues equal to 1 corresponding to real eigenvectors that span  $v^\perp$ . The remaining eigenvalue

$$\lambda = 1 - \zeta v^T u \quad (1796)$$

corresponds to eigenvector  $u$ .<sup>B.8</sup> From [248, App.7.A.26] the determinant:

$$\det E = 1 - \text{tr}(\zeta uv^T) = \lambda \quad (1797)$$

If  $\lambda \neq 0$  then  $E$  is invertible; [174] (confer §B.1.0.1)

$$E^{-1} = I + \frac{\zeta}{\lambda} uv^T \quad (1798)$$

Eigenvectors corresponding to 0 eigenvalues belong to  $\mathcal{N}(E)$ , and the number of 0 eigenvalues must be at least  $\dim \mathcal{N}(E)$  which, here, can be at most one.

<sup>B.8</sup>Elementary matrix  $E$  is not always diagonalizable because eigenvector  $u$  need not be independent of the others; *id est*,  $u \in v^\perp$  is possible.

(§A.7.3.0.1) The nullspace exists, therefore, when  $\lambda=0$ ; *id est*, when  $v^T u = 1/\zeta$ ; rather, whenever  $u$  belongs to hyperplane  $\{z \in \mathbb{R}^N \mid v^T z = 1/\zeta\}$ . Then (when  $\lambda=0$ ) elementary matrix  $E$  is a nonorthogonal projector projecting on its range ( $E^2=E$ , §E.1) and  $\mathcal{N}(E)=\mathcal{R}(u)$ ; eigenvector  $u$  spans the nullspace when it exists. By conservation of dimension,  $\dim \mathcal{R}(E)=N-\dim \mathcal{N}(E)$ . It is apparent from (1795) that  $v^\perp \subseteq \mathcal{R}(E)$ , but  $\dim v^\perp=N-1$ . Hence  $\mathcal{R}(E) \equiv v^\perp$  when the nullspace exists, and the remaining eigenvectors span it.

In summary, when a nontrivial nullspace of  $E$  exists,

$$\mathcal{R}(E) = \mathcal{N}(v^T), \quad \mathcal{N}(E) = \mathcal{R}(u), \quad v^T u = 1/\zeta \quad (1799)$$

illustrated in Figure 183, which is opposite to the assignment of subspaces for a dyad (Figure 181). Otherwise,  $\mathcal{R}(E)=\mathbb{R}^N$ .

When  $E=E^T$ , the spectral norm is

$$\|E\|_2 = \max\{1, |\lambda|\} \quad (1800)$$

### B.3.1 Householder matrix

An elementary matrix is called a Householder matrix when it has the defining form, for nonzero vector  $u$  [181, §5.1.2] [174, §4.10.1] [368, §7.3] [228, §2.2]

$$H = I - 2 \frac{uu^T}{u^T u} \in \mathbb{S}^N \quad (1801)$$

which is a symmetric orthogonal (reflection) matrix ( $H^{-1}=H^T=H$  (§B.5.3)). Vector  $u$  is normal to an  $N-1$ -dimensional subspace  $u^\perp$  through which this particular  $H$  effects pointwise reflection; *e.g.*,  $Hu^\perp=u^\perp$  while  $Hu=-u$ .

Matrix  $H$  has  $N-1$  orthonormal eigenvectors spanning that reflecting subspace  $u^\perp$  with corresponding eigenvalues equal to 1. The remaining eigenvector  $u$  has corresponding eigenvalue  $-1$ ; so

$$\det H = -1 \quad (1802)$$

Due to symmetry of  $H$ , the matrix 2-norm (the spectral norm) is equal to the largest eigenvalue-magnitude. A Householder matrix is thus characterized,

$$H^T = H, \quad H^{-1} = H^T, \quad \|H\|_2 = 1, \quad H \neq 0 \quad (1803)$$

For example, the permutation matrix

$$\Xi = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad (1804)$$

is a Householder matrix having  $u=[0 \ 1 \ -1]^T/\sqrt{2}$ . Not all permutation matrices are Householder matrices, although all permutation matrices are orthogonal matrices (§B.5.2,  $\Xi^T \Xi = I$ ) [368, §3.4] because they are made by permuting rows and columns of the Identity matrix. Neither are all symmetric permutation matrices Householder matrices;

*e.g.*,  $\Xi = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$  (1900) is not a Householder matrix.

## B.4 Auxiliary $V$ -matrices

### B.4.1 Auxiliary projector matrix $V$

It is convenient to define a matrix  $V$  that arises naturally as a consequence of translating geometric center  $\alpha_c$  (§5.5.1.0.1) of some list  $X$  to the origin. In place of  $X - \alpha_c \mathbf{1}^T$  we may write  $XV$  as in (1119) where

$$V = I - \frac{1}{N} \mathbf{1} \mathbf{1}^T \in \mathbb{S}^N \quad (1055)$$

is an elementary matrix called the *geometric centering matrix*.

Any elementary matrix in  $\mathbb{R}^{N \times N}$  has  $N-1$  eigenvalues equal to 1. For the particular elementary matrix  $V$ , the  $N^{\text{th}}$  eigenvalue equals 0. The number of 0 eigenvalues must equal  $\dim \mathcal{N}(V) = 1$ , by the 0 *eigenvalues theorem* (§A.7.3.0.1), because  $V = V^T$  is diagonalizable. Because

$$V \mathbf{1} = \mathbf{0} \quad (1805)$$

the nullspace  $\mathcal{N}(V) = \mathcal{R}(\mathbf{1})$  is spanned by the eigenvector  $\mathbf{1}$ . The remaining eigenvectors span  $\mathcal{R}(V) \equiv \mathbf{1}^\perp = \mathcal{N}(\mathbf{1}^T)$  that has dimension  $N-1$ .

Because

$$V^2 = V \quad (1806)$$

and  $V^T = V$ , elementary matrix  $V$  is also a projection matrix (§E.3) projecting orthogonally on its range  $\mathcal{N}(\mathbf{1}^T)$  which is a hyperplane containing the origin in  $\mathbb{R}^N$

$$V = I - \mathbf{1}(\mathbf{1}^T \mathbf{1})^{-1} \mathbf{1}^T \quad (1807)$$

The  $\{0, 1\}$  eigenvalues also indicate that diagonalizable  $V$  is a projection matrix. [455, §4.1 thm.4.1] Symmetry of  $V$  denotes orthogonal projection; from (2100),

$$V^2 = V, \quad V^T = V, \quad V^\dagger = V, \quad \|V\|_2 = 1, \quad V \succeq 0 \quad (1808)$$

Matrix  $V$  is also circulant [192].

#### B.4.1.0.1 Example. Relationship of Auxiliary to Householder matrix.

Let  $H \in \mathbb{S}^N$  be a Householder matrix (1801) defined by

$$u = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ 1 + \sqrt{N} \end{bmatrix} \in \mathbb{R}^N \quad (1809)$$

Then we have [176, §2]

$$V = H \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0}^T & 0 \end{bmatrix} H \quad (1810)$$

Let  $D \in \mathbb{S}_h^N$  and define

$$-HDH \triangleq - \begin{bmatrix} A & b \\ b^T & c \end{bmatrix} \quad (1811)$$

where  $b$  is a vector. Then because  $H$  is nonsingular (§A.3.1.0.5) [210, §3]

$$-VDV = -H \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0}^T & 0 \end{bmatrix} H \succeq 0 \Leftrightarrow -A \succeq 0 \quad (1812)$$

and affine dimension is  $r = \text{rank } A$  when  $D$  is a Euclidean distance matrix.  $\square$

### B.4.2 Schoenberg auxiliary matrix $V_{\mathcal{N}}$

1.  $V_{\mathcal{N}} = \frac{1}{\sqrt{2}} \begin{bmatrix} -\mathbf{1}^T \\ I \end{bmatrix} \in \mathbb{R}^{N \times N-1}$  (1039)
2.  $V_{\mathcal{N}}^T \mathbf{1} = \mathbf{0}$
3.  $I - e_1 \mathbf{1}^T = [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]$
4.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] V_{\mathcal{N}} = V_{\mathcal{N}}$
5.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] V = V$
6.  $V [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] = [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]$
7.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] = [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]$
8.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]^{\dagger} = \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} V$
9.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]^{\dagger} V = [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]^{\dagger}$
10.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]^{\dagger} = V$
11.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]^{\dagger} [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] = \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix}$
12.  $[\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} = [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}]$
13.  $\begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} [\mathbf{0} \quad \sqrt{2}V_{\mathcal{N}}] = \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix}$
14.  $[V_{\mathcal{N}} \quad \frac{1}{\sqrt{2}}\mathbf{1}]^{-1} = \begin{bmatrix} V_{\mathcal{N}}^{\dagger} \\ \frac{\sqrt{2}}{N}\mathbf{1}^T \end{bmatrix}$
15.  $V_{\mathcal{N}}^{\dagger} = \sqrt{2} [-\frac{1}{N}\mathbf{1} \quad I - \frac{1}{N}\mathbf{1}\mathbf{1}^T] \in \mathbb{R}^{N-1 \times N}, \quad \left( I - \frac{1}{N}\mathbf{1}\mathbf{1}^T \in \mathbb{S}^{N-1} \right)$
16.  $V_{\mathcal{N}}^{\dagger} \mathbf{1} = \mathbf{0}$
17.  $V_{\mathcal{N}}^{\dagger} V_{\mathcal{N}} = I, \quad V_{\mathcal{N}}^T V_{\mathcal{N}} = \frac{1}{2}(I + \mathbf{1}\mathbf{1}^T) \in \mathbb{S}^{N-1}$
18.  $V^T = V = V_{\mathcal{N}} V_{\mathcal{N}}^{\dagger} = I - \frac{1}{N}\mathbf{1}\mathbf{1}^T \in \mathbb{S}^N$
19.  $-V_{\mathcal{N}}^{\dagger}(\mathbf{1}\mathbf{1}^T - I)V_{\mathcal{N}} = I, \quad (\mathbf{1}\mathbf{1}^T - I \in \mathbb{EDM}^N)$
20.  $D = [d_{ij}] \in \mathbb{S}_h^N \quad (1057)$   
 $\text{tr}(-VDV) = \text{tr}(-VD) = \text{tr}(-V_{\mathcal{N}}^{\dagger}DV_{\mathcal{N}}) = \frac{1}{N}\mathbf{1}^T D \mathbf{1} = \frac{1}{N} \text{tr}(\mathbf{1}\mathbf{1}^T D) = \frac{1}{N} \sum_{i,j} d_{ij}$

Any elementary matrix  $E \in \mathbb{S}^N$  of the particular form

$$E = k_1 I - k_2 \mathbf{1}\mathbf{1}^T \quad (1813)$$

where  $k_1, k_2 \in \mathbb{R}$ , [B.9](#) will make  $\text{tr}(-ED)$  proportional to  $\sum d_{ij}$ .

---

[B.9](#)If  $k_1$  is  $1-\rho$  while  $k_2$  equals  $-\rho \in \mathbb{R}$ , then all eigenvalues of  $E$  for  $-1/(N-1) < \rho < 1$  are guaranteed positive and therefore  $E$  is guaranteed positive definite. [\[337\]](#)

$$21. \quad D = [d_{ij}] \in \mathbb{S}^N$$

$$\text{tr}(-VDV) = \frac{1}{N} \sum_{\substack{i,j \\ i \neq j}} d_{ij} - \frac{N-1}{N} \sum_i d_{ii} = \frac{1}{N} \mathbf{1}^T D \mathbf{1} - \text{tr} D$$

$$22. \quad D = [d_{ij}] \in \mathbb{S}_h^N$$

$$\text{tr}(-V_N^T D V_N) = \sum_j d_{1j}$$

$$23. \text{ For } Y \in \mathbb{S}^N$$

$$V(Y - \delta(Y\mathbf{1}))V = Y - \delta(Y\mathbf{1})$$

### B.4.3 Orthonormal auxiliary matrix $V_W$

Thin matrix

$$V_W \triangleq \begin{bmatrix} \frac{-1}{\sqrt{N}} & \frac{-1}{\sqrt{N}} & \cdots & \frac{-1}{\sqrt{N}} \\ 1 + \frac{-1}{N+\sqrt{N}} & \frac{-1}{N+\sqrt{N}} & \cdots & \frac{-1}{N+\sqrt{N}} \\ \frac{-1}{N+\sqrt{N}} & \ddots & \ddots & \frac{-1}{N+\sqrt{N}} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{-1}{N+\sqrt{N}} & \frac{-1}{N+\sqrt{N}} & \cdots & 1 + \frac{-1}{N+\sqrt{N}} \end{bmatrix} \in \mathbb{R}^{N \times N-1} \quad (1814)$$

has  $\mathcal{R}(V_W) = \mathcal{N}(\mathbf{1}^T)$  and orthonormal columns. [7] We defined three auxiliary  $V$ -matrices:  $V$ ,  $V_N$  (1039), and  $V_W$  sharing some attributes listed in Table B.4.4. For example,  $V$  can be expressed

$$V = V_W V_W^T = V_N V_N^\dagger \quad (1815)$$

but  $V_W^T V_W = I$  means  $V$  is an orthogonal projector (2097) and

$$V_W^\dagger = V_W^T, \quad \|V_W\|_2 = 1, \quad V_W^T \mathbf{1} = \mathbf{0} \quad (1816)$$

### B.4.4 Auxiliary $V$ -matrix Table

	$\dim V$	$\text{rank } V$	$\mathcal{R}(V)$	$\mathcal{N}(V^T)$	$V^T V$	$V V^T$	$V V^\dagger$
$V$	$N \times N$	$N-1$	$\mathcal{N}(\mathbf{1}^T)$	$\mathcal{R}(\mathbf{1})$	$V$	$V$	$V$
$V_N$	$N \times (N-1)$	$N-1$	$\mathcal{N}(\mathbf{1}^T)$	$\mathcal{R}(\mathbf{1})$	$\frac{1}{2}(I + \mathbf{1}\mathbf{1}^T)$	$\frac{1}{2} \begin{bmatrix} N-1 & -\mathbf{1}^T \\ -\mathbf{1} & I \end{bmatrix}$	$V$
$V_W$	$N \times (N-1)$	$N-1$	$\mathcal{N}(\mathbf{1}^T)$	$\mathcal{R}(\mathbf{1})$	$I$	$V$	$V$

### B.4.5 More auxiliary matrices

Mathar shows [290, §2] that any elementary matrix (§B.3) of the form

$$V_S = I - b \mathbf{1}^T \in \mathbb{R}^{N \times N} \quad (1817)$$

such that  $b^T \mathbf{1} = 1$  (confer [186, §2]), is an auxiliary  $V$ -matrix having

$$\begin{aligned} \mathcal{R}(V_S^T) &= \mathcal{N}(b^T), & \mathcal{R}(V_S) &= \mathcal{N}(\mathbf{1}^T) \\ \mathcal{N}(V_S) &= \mathcal{R}(b), & \mathcal{N}(V_S^T) &= \mathcal{R}(\mathbf{1}) \end{aligned} \quad (1818)$$

Given  $X \in \mathbb{R}^{n \times N}$ , the choice  $b = \frac{1}{N}\mathbf{1}$  ( $V_S = V$ ) minimizes  $\|X(I - b\mathbf{1}^T)\|_F$ . [188, §3.2.1]

## B.5 Orthomatrices

### B.5.1 Orthonormal matrix

Property  $Q^T Q = I$  completely defines orthonormal matrix  $Q \in \mathbb{R}^{n \times k}$  ( $k \leq n$ ); a thin-or-square full-rank matrix characterized by nonexpansivity (2101)

$$\|Q^T x\|_2 \leq \|x\|_2 \quad \forall x \in \mathbb{R}^n, \quad \|Qy\|_2 = \|y\|_2 \quad \forall y \in \mathbb{R}^k \quad (1819)$$

and preservation of vector inner-product

$$\langle Qy, Qz \rangle = \langle y, z \rangle \quad (1820)$$

### B.5.2 Orthogonal matrix & vector rotation

An orthogonal matrix is a square orthonormal matrix. Property  $Q^{-1} = Q^T$  completely defines orthogonal matrix  $Q \in \mathbb{R}^{n \times n}$  employed to effect vector rotation; [368, §2.6, §3.4] [370, §6.5] [228, §2.1] for any  $x \in \mathbb{R}^n$

$$\|Qx\|_2 = \|x\|_2 \quad (1821)$$

In other words, the 2-norm is orthogonally invariant. Any antisymmetric matrix constructs an orthogonal matrix; *id est*, for  $A = -A^T$

$$Q = (I + A)^{-1}(I - A) \quad (1822)$$

A *unitary matrix* is a complex generalization of orthogonal matrix; conjugate transpose defines it:  $U^{-1} = U^H$ . An orthogonal matrix is simply a real unitary matrix. **B.10**

Orthogonal matrix  $Q$  is a normal matrix further characterized by spectral norm:

$$Q^{-1} = Q^T, \quad \|Q\|_2 = 1 \quad (1823)$$

Applying characterization (1823) to  $Q^T$ , we see it too is an orthogonal matrix. Hence the rows and columns of  $Q$  respectively form an orthonormal set. Normalcy guarantees diagonalization (§A.5.1) so, for  $Q \triangleq S\Lambda S^H$

$$S\Lambda^{-1}S^H = S^*\Lambda S^T (= S\Lambda^*S^H), \quad \|\delta(\Lambda)\|_\infty = 1 \quad (1824)$$

characterizes an orthogonal matrix in terms of eigenvalues and eigenvectors.

All permutation matrices  $\Xi$ , for example, are nonnegative orthogonal matrices; and *vice versa*. Product or Kronecker product of any permutation matrices remains a permutator. Any product of permutation matrix with orthogonal matrix remains orthogonal. In fact, any product  $AQ$  of orthogonal matrices  $A$  and  $Q$  remains orthogonal by definition. Given any other dimensionally compatible orthogonal matrix  $U$ , the mapping  $g(A) = U^T AQ$  is a bijection on the domain of orthogonal matrices (a nonconvex manifold of dimension  $\frac{1}{2}n(n-1)$  [55]). [268, §2.1] [269]

The largest magnitude entry of an orthogonal matrix is 1; for each and every  $j \in 1 \dots n$

$$\begin{aligned} \|Q(j, :)\|_\infty &\leq 1 \\ \|Q(:, j)\|_\infty &\leq 1 \end{aligned} \quad (1825)$$

Each and every eigenvalue of a (real) orthogonal matrix has magnitude 1 ( $\Lambda^{-1} = \Lambda^*$ )

$$\lambda(Q) \in \mathbb{C}^n, \quad |\lambda(Q)| = 1 \quad (1826)$$

but only the Identity matrix can be simultaneously orthogonal and positive definite. Orthogonal matrices have complex eigenvalues in conjugate pairs: so  $\det Q = \pm 1$ .

---

**B.10**Orthogonal and unitary matrices are called *unitary linear operators*.

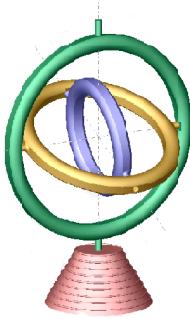


Figure 184: *Gimbal*: a mechanism imparting three degrees of dimensional freedom to a Euclidean body suspended at its center. Each ring is free to rotate about one axis. (Drawing by courtesy of [The MathWorks Inc.](#))

### B.5.3 Reflection

A matrix for pointwise reflection is defined by imposing symmetry upon the orthogonal matrix; *id est*, a *reflection matrix* is completely defined by  $Q^{-1} = Q^T = Q$ . The reflection matrix is a symmetric orthogonal matrix, and *vice versa*, characterized:

$$Q^T = Q, \quad Q^{-1} = Q^T, \quad \|Q\|_2 = 1 \quad (1827)$$

The Householder matrix (§B.3.1) is an example of symmetric orthogonal (reflection) matrix.

Reflection matrices have eigenvalues equal to  $\pm 1$  and so  $\det Q = \pm 1$ . It is natural to expect a relationship between reflection and projection matrices because all projection matrices have eigenvalues belonging to  $\{0, 1\}$ . In fact, any reflection matrix  $Q$  is related to some orthogonal projector  $P$  by [230, §1 prob.44]

$$Q = I - 2P \quad (1828)$$

Yet  $P$  is, generally, neither orthogonal or invertible. (§E.3.2)

$$\lambda(Q) \in \mathbb{R}^n, \quad |\lambda(Q)| = 1 \quad (1829)$$

Reflection is with respect to  $\mathcal{R}(P)^\perp$ . Matrix  $2P - I$  represents antireflection.

Every orthogonal matrix can be expressed as the product of a rotation and a reflection. The collection of all orthogonal matrices of particular dimension does not form a convex set.

### B.5.4 Rotation of range and rowspace

Given orthogonal matrix  $Q$ , column vectors of a matrix  $X$  are simultaneously rotated about the origin via product  $QX$ . In three dimensions ( $X \in \mathbb{R}^{3 \times N}$ ), the precise meaning of rotation is best illustrated in Figure 184 where a gimbal aids visualization of what is achievable; mathematically, (§5.5.2.0.1)

$$Q = \begin{bmatrix} \cos \theta & 0 & -\sin \theta \\ 0 & 1 & 0 \\ \sin \theta & 0 & \cos \theta \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \psi & -\sin \psi \\ 0 & \sin \psi & \cos \psi \end{bmatrix} \begin{bmatrix} \cos \phi & -\sin \phi & 0 \\ \sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (1830)$$

#### B.5.4.0.1 Example. One axis of revolution.

Partition  $n+1$ -dimensional Euclidean space  $\mathbb{R}^{n+1} \triangleq \begin{bmatrix} \mathbb{R}^n \\ \mathbb{R} \end{bmatrix}$  and define an  $n$ -dimensional subspace

$$\mathcal{R} \triangleq \{\lambda \in \mathbb{R}^{n+1} \mid \mathbf{1}^T \lambda = 0\} \quad (1831)$$

(a hyperplane through the origin). We want an orthogonal matrix that rotates a list in the columns of matrix  $X \in \mathbb{R}^{n+1 \times N}$  through the dihedral angle between  $\mathbb{R}^n$  and  $\mathcal{R}$  (§2.4.3)

$$\alpha(\mathbb{R}^n, \mathcal{R}) = \arccos\left(\frac{\langle e_{n+1}, \mathbf{1} \rangle}{\|e_{n+1}\| \|\mathbf{1}\|}\right) = \arccos\left(\frac{1}{\sqrt{n+1}}\right) \text{ radians} \quad (1832)$$

The vertex-description of the nonnegative orthant in  $\mathbb{R}^{n+1}$  is

$$\{[e_1 \ e_2 \ \cdots \ e_{n+1}]a \mid a \succeq 0\} = \{a \succeq 0\} = \mathbb{R}_+^{n+1} \subset \mathbb{R}^{n+1} \quad (1833)$$

Consider rotation of these vertices via orthogonal matrix

$$Q \triangleq [\mathbf{1} \frac{1}{\sqrt{n+1}} \ \Xi V_{\mathcal{W}}] \Xi \in \mathbb{R}^{n+1 \times n+1} \quad (1834)$$

where permutation matrix  $\Xi \in \mathbb{S}^{n+1}$  is defined in (1900), and  $V_{\mathcal{W}} \in \mathbb{R}^{n+1 \times n}$  is the orthonormal auxiliary matrix defined in §B.4.3. This particular orthogonal matrix is selected because it rotates any point in subspace  $\mathbb{R}^n$  about one axis of revolution onto  $\mathcal{R}$ ; e.g., rotation  $Qe_{n+1}$  aligns the last standard basis vector with subspace normal  $\mathcal{R}^\perp = \mathbf{1}$ . The rotated standard basis vectors remaining are orthonormal spanning  $\mathcal{R}$ .  $\square$

Another interpretation of product  $QX$  is rotation/reflection of  $\mathcal{R}(X)$ . Rotation of  $X$  as in  $QXQ^T$  is a simultaneous rotation/reflection of range and rowspace. [B.11](#)

**Proof.** Any matrix can be expressed as a singular value decomposition  $X = U\Sigma W^T$  (1716) where  $\delta^2(\Sigma) = \Sigma$ ,  $\mathcal{R}(U) \supseteq \mathcal{R}(X)$ , and  $\mathcal{R}(W) \supseteq \mathcal{R}(X^T)$ .  $\spadesuit$

#### B.5.5 Matrix rotation

Orthogonal matrices are also employed to rotate/reflect other matrices like vectors: [181, §12.4.1] Given orthogonal matrix  $Q$ , the product  $Q^TA$  will rotate  $A \in \mathbb{R}^{n \times n}$  in the Euclidean sense in  $\mathbb{R}^{n^2}$  because Frobenius' norm is orthogonally invariant (§2.2.1);

$$\|Q^TA\|_F = \sqrt{\text{tr}(A^TQQ^TA)} = \|A\|_F \quad (1835)$$

(likewise for  $AQ$ ). Were  $A$  symmetric, such a rotation would depart from  $\mathbb{S}^n$ . One remedy is to instead form product  $Q^TAQ$  because

$$\|Q^TAQ\|_F = \sqrt{\text{tr}(Q^TA^TQQ^TAQ)} = \|A\|_F \quad (1836)$$

By [§A.1.1 no.33](#),

$$\text{vec } Q^TAQ = (Q \otimes Q)^T \text{vec } A \quad (1837)$$

which is a rotation of the vectorized  $A$  matrix because Kronecker product of any orthogonal matrices remains orthogonal; e.g., by [§A.1.1 no.43](#),

$$(Q \otimes Q)^T(Q \otimes Q) = I \quad (1838)$$

Matrix  $A$  is *orthogonally equivalent* to  $B$  if  $B = S^TAS$  for some orthogonal matrix  $S$ . Every square matrix, for example, is orthogonally equivalent to a matrix having equal entries along the main diagonal. [228, §2.2, prob.3]

---

[B.11](#) Product  $Q^TAQ$  can be regarded as coordinate transformation; e.g., given linear map  $y = Ax : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and orthogonal  $Q$ , the transformation  $Qy = AQx$  is a rotation/reflection of range and rowspace (143) of matrix  $A$  where  $Qy \in \mathcal{R}(A)$  and  $Qx \in \mathcal{R}(A^T)$  (144).

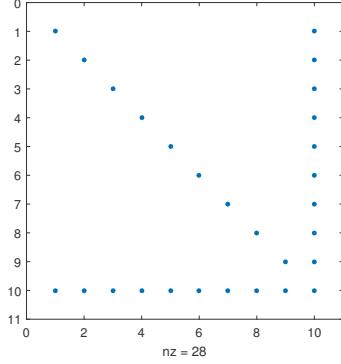


Figure 185:  $10 \times 10$  arrow matrix. Twenty eight nonzero ( $\text{nz}$ ) entries indicated.

## B.6 Arrow matrix

Consider a partitioned symmetric  $n \times n$ -dimensional matrix  $A$  that has *arrow* [324] (or *arrowhead* [367]) form constituted by vectors  $a, b \in \mathbb{R}^{n-1}$  and real scalar  $c$ :

$$A \triangleq \begin{bmatrix} \delta(a) & b \\ b^T & c \end{bmatrix} \in \mathbb{S}^n \quad (1839)$$

Figure 185 illustrates sparsity pattern of an arrow matrix. Embedding of diagonal matrix  $\delta(a)$  makes relative sparsity increasing with dimension. Because an arrow matrix is a kind of bordered matrix, eigenvalues of  $\delta(a)$  and  $A$  are interlaced;

$$\lambda_n \leq (\Xi^T a)_{n-1} \leq \lambda_{n-1} \leq (\Xi^T a)_{n-2} \leq \cdots \leq (\Xi^T a)_1 \leq \lambda_1 \quad (1840)$$

[368, §6.4] [228, §4.3] [364, §IV.4.1] denoting nonincreasingly ordered eigenvalues of  $A$  by vector  $\lambda \in \mathbb{R}^n$ , and those of  $\delta(a)$  by  $\Xi^T a \in \mathbb{R}^{n-1}$  where  $\Xi$  is a permutation matrix arranging  $a$  into nonincreasing order:  $\delta(a) = \Xi \delta(\Xi^T a) \Xi^T$  (§A.5.1.2).

### B.6.1 positive semidefinite arrow matrix

i) Nonnegative main diagonal  $a \succeq 0$  insures  $n-1$  nonnegative eigenvalues in (1840).

Positive semidefiniteness is left determined by smallest eigenvalue  $\lambda_n$ :

$$\begin{aligned} A \succeq 0 &\Leftrightarrow a \succeq 0, \quad b^T(I - \delta(a)\delta(a)^\dagger) = \mathbf{0}, \quad c - b^T\delta(a)^\dagger b \geq 0 \\ &\Leftrightarrow c \geq 0, \quad b(1 - cc^\dagger) = \mathbf{0}, \quad \delta(a) - c^\dagger bb^T \succeq 0 \end{aligned} \quad (1841)$$

Schur complement condition (§A.4)  $b^T(I - \delta(a)\delta(a)^\dagger) = \mathbf{0}$  is most simply a requirement for

ii) a zero entry in vector  $b$  wherever there is a corresponding zero entry in vector  $a$ .

In other words, vector  $b$  can reside anywhere in a Cartesian subspace of  $\mathbb{R}^{n-1}$  that is determined solely by indices of the nonzero entries in vector  $a$ .

iii)  $c \geq b^T\delta(a)^\dagger b$  provides a tight lower bound for scalar  $c$ .

As shown in §3.5.1,  $b^T\delta(a)^\dagger b$  is simultaneously convex in vectors  $a$  and  $b$ .

## Appendix C

# Some analytical optimal results

*People have been working on Optimization since the ancient Greeks [Zenodorus, circa 200BC] learned that a string encloses the most area when it is formed into the shape of a circle.*

—ROMAN POLYAK

We speculate that optimization problems possessing analytical solution have convex transformation or constructive global optimality conditions, perhaps yet unknown; e.g, §4.10.2, §7.1.4, (1872), §C.3.0.1.

### C.1 Properties of infima

- 

$$\begin{aligned} \inf \emptyset &\triangleq \infty \\ \sup \emptyset &\triangleq -\infty \end{aligned} \tag{1842}$$

- Given  $f(x) : \mathcal{X} \rightarrow \mathbb{R}$  defined on arbitrary set  $\mathcal{X}$  [225, §0.1.2]

$$\begin{aligned} \inf_{x \in \mathcal{X}} f(x) &= -\sup_{x \in \mathcal{X}} -f(x) \\ \sup_{x \in \mathcal{X}} f(x) &= -\inf_{x \in \mathcal{X}} -f(x) \end{aligned} \tag{1843}$$

$$\begin{aligned} \arg \inf_{x \in \mathcal{X}} f(x) &= \arg \sup_{x \in \mathcal{X}} -f(x) \\ \arg \sup_{x \in \mathcal{X}} f(x) &= \arg \inf_{x \in \mathcal{X}} -f(x) \end{aligned} \tag{1844}$$

- Given scalar  $\kappa$  and  $f(x) : \mathcal{X} \rightarrow \mathbb{R}$  and  $g(x) : \mathcal{X} \rightarrow \mathbb{R}$  defined on arbitrary set  $\mathcal{X}$  [225, §0.1.2]

$$\begin{aligned} \inf_{x \in \mathcal{X}} (\kappa + f(x)) &= \kappa + \inf_{x \in \mathcal{X}} f(x) \\ \arg \inf_{x \in \mathcal{X}} (\kappa + f(x)) &= \arg \inf_{x \in \mathcal{X}} f(x) \end{aligned} \tag{1845}$$

$$\left. \begin{aligned} \inf_{x \in \mathcal{X}} \kappa f(x) &= \kappa \inf_{x \in \mathcal{X}} f(x) \\ \arg \inf_{x \in \mathcal{X}} \kappa f(x) &= \arg \inf_{x \in \mathcal{X}} f(x) \end{aligned} \right\}, \quad \kappa > 0 \tag{1846}$$

$$\inf_{x \in \mathcal{X}} (f(x) + g(x)) \geq \inf_{x \in \mathcal{X}} f(x) + \inf_{x \in \mathcal{X}} g(x) \tag{1847}$$

- Given  $f(x) : \mathcal{X} \rightarrow \mathbb{R}$  defined on arbitrary set  $\mathcal{X}$

$$\arg \inf_{x \in \mathcal{X}} |f(x)| = \arg \inf_{x \in \mathcal{X}} f(x)^2 \quad (1848)$$

- Given  $f(x) : \mathcal{X} \cup \mathcal{Y} \rightarrow \mathbb{R}$  and arbitrary sets  $\mathcal{X}$  and  $\mathcal{Y}$  [225, §0.1.2]

$$\mathcal{X} \subset \mathcal{Y} \Rightarrow \inf_{x \in \mathcal{X}} f(x) \geq \inf_{x \in \mathcal{Y}} f(x) \quad (1849)$$

$$\inf_{x \in \mathcal{X} \cup \mathcal{Y}} f(x) = \min\{\inf_{x \in \mathcal{X}} f(x), \inf_{x \in \mathcal{Y}} f(x)\} \quad (1850)$$

$$\inf_{x \in \mathcal{X} \cap \mathcal{Y}} f(x) \geq \max\{\inf_{x \in \mathcal{X}} f(x), \inf_{x \in \mathcal{Y}} f(x)\} \quad (1851)$$

## C.2 Trace, singular and eigen values

- For  $A \in \mathbb{R}^{m \times n}$  and  $\sigma(A)$  denoting its singular values, the *nuclear* (Ky Fan) *norm*  $\|A\|_2^*$  of matrix  $A$  (confer (44), (1719), [229, p.200]) is

$$\begin{aligned} \sum_i \sigma(A)_i &= \text{tr} \sqrt{A^T A} = \|A\|_2^* = \sup_{\|X\|_2 \leq 1} \text{tr}(X^T A) = \underset{X \in \mathbb{R}^{m \times n}}{\text{maximize}} \quad \text{tr}(X^T A) \\ &\quad \text{subject to } \begin{bmatrix} I & X \\ X^T & I \end{bmatrix} \succeq 0 \\ &= \frac{1}{2} \underset{X \in \mathbb{S}^m, Y \in \mathbb{S}^n}{\text{minimize}} \quad \text{tr } X + \text{tr } Y \\ &\quad \text{subject to } \begin{bmatrix} X & A \\ A^T & Y \end{bmatrix} \succeq 0 \end{aligned} \quad (1852)$$

This nuclear norm is convex<sup>C.1</sup> and dual to the spectral norm. [229, p.214] [65, §A.1.6] Given singular value decomposition  $A = S\Sigma Q^T \in \mathbb{R}^{m \times n}$  (A.6), then  $X^* = SQ^T \in \mathbb{R}^{m \times n}$  is an optimal solution to maximization (confer §2.3.2.0.5) while  $X^* = S\Sigma S^T \in \mathbb{S}^m$  and  $Y^* = Q\Sigma Q^T \in \mathbb{S}^n$  is an optimal solution to minimization [153]. Srebro [359] asserts

$$\begin{aligned} \sum_i \sigma(A)_i &= \frac{1}{2} \underset{U, V}{\text{minimize}} \quad \|U\|_F^2 + \|V\|_F^2 \\ &\quad \text{subject to } A = UV^T \\ &= \underset{U, V}{\text{minimize}} \quad \|U\|_F \|V\|_F \\ &\quad \text{subject to } A = UV^T \end{aligned} \quad (1853)$$

- For  $A \in \mathbb{R}^{m \times n}$  and  $\sigma(A)_1$  connoting spectral norm,

$$\sigma(A)_1 = \sqrt{\lambda(A^T A)_1} = \|A\|_2 = \sup_{\|x\|=1} \|Ax\|_2 = \underset{t \in \mathbb{R}}{\text{minimize}} \quad t \quad (596)$$

$$\text{subject to } \begin{bmatrix} tI & A \\ A^T & tI \end{bmatrix} \succeq 0$$

Denoting  $\rho = \text{rank } A$

$$\sigma(A)_\rho = \sqrt{\lambda(A^T A)_\rho} = \|A^\dagger\|_2^{-1} = 1 / \underset{t \in \mathbb{R}}{\text{minimize}} \quad t \quad (1854)$$

$$\text{subject to } \begin{bmatrix} tI & A^\dagger \\ A^{\dagger T} & tI \end{bmatrix} \succeq 0$$

which is equal to  $\inf_{\|x\|=1} \|Ax\|_2$  when  $A$  is full rank; *id est*, when  $\rho = \min\{m, n\}$ .

---

<sup>C.1</sup> discernible as envelope of the rank function (1516) or as supremum of functions linear in  $A$  (Figure 78).

By confining dyad  $uv^T$  to the unit nuclear norm ball (95),

$$\begin{aligned} \sigma(A)_1 = \|A\|_2 &= \sup_{\|u\|=1, \|v\|=1} u^T A v = \underset{\substack{Z \in \mathbb{R}^{m \times n}, X \in \mathbb{S}^m, Y \in \mathbb{S}^n \\ \text{subject to}}} {\underset{\text{maximize}}{\text{tr}}}(Z^T A) \\ &\quad \text{subject to} \quad \text{tr } X + \text{tr } Y \leq 2 \\ &\quad \left[ \begin{array}{cc} X & Z \\ Z^T & Y \end{array} \right] \succeq 0 \end{aligned} \quad (1855)$$

with corresponding left and right singular vectors (optimal)  $u^*$  and  $v^*$ . Applying (1862) to a result of Lanczos [180, p.207],

$$\begin{aligned} \sigma(A)_1 = \|A\|_2 &= \sup_{\left\| \begin{bmatrix} u \\ v \end{bmatrix} \right\|=1} \begin{bmatrix} u \\ v \end{bmatrix}^T \begin{bmatrix} \mathbf{0} & A \\ A^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \underset{\substack{X \in \mathbb{S}_+^{m+n} \\ \text{subject to}}} {\underset{\text{maximize}}{\text{tr}}}\left(X \begin{bmatrix} \mathbf{0} & A \\ A^T & \mathbf{0} \end{bmatrix}\right) \\ &\quad \text{subject to} \quad \text{tr } X = 1 \\ &= \underset{t \in \mathbb{R}}{\underset{\substack{\text{subject to}}} {\underset{\text{minimize}}{\text{tr}}}} t \\ &\quad \left[ \begin{array}{cc} \mathbf{0} & A \\ A^T & \mathbf{0} \end{array} \right] \preceq t I \end{aligned} \quad (1856)$$

whose corresponding left and right singular vectors are  $\sqrt{2}u^*$  and  $\sqrt{2}v^*$ .

#### C.2.0.0.1 Exercise. Optimal matrix factorization.

Prove (1853). ▼

- For  $X \in \mathbb{S}^m$ ,  $Y \in \mathbb{S}^n$ ,  $A \in \mathcal{C} \subseteq \mathbb{R}^{m \times n}$  for set  $\mathcal{C}$  convex, and  $\sigma(A)$  denoting the singular values of  $A$  [153, §3],

$$\begin{aligned} \underset{A}{\underset{\substack{\text{minimize} \\ \text{subject to}}} {\text{minimize}}} \sum_i \sigma(A)_i &\equiv \underset{\substack{A, X, Y \\ A \in \mathcal{C}}}{\underset{\substack{\frac{1}{2} \text{ minimize} \\ \text{subject to}}} {\text{minimize}}} \text{tr } X + \text{tr } Y \\ &\quad \left[ \begin{array}{cc} X & A \\ A^T & Y \end{array} \right] \succeq 0 \end{aligned} \quad (1857)$$

For feasible set  $\mathcal{C}$  equal to the unit nuclear norm ball (95),

$$\begin{aligned} \text{find } A &\quad \underset{\substack{A, X, Y \\ \text{subject to}}} {\underset{\substack{\text{find} \\ A \in \{Z \in \mathbb{R}^{m \times n} \mid \sum_i \sigma(Z)_i \leq 1\}}} {\text{find}}} A \\ \text{subject to } A \in \{Z \in \mathbb{R}^{m \times n} \mid \sum_i \sigma(Z)_i \leq 1\} &\equiv \underset{\substack{\text{subject to} \\ \left[ \begin{array}{cc} X & A \\ A^T & Y \end{array} \right] \succeq 0}} {\underset{\text{subject to}} {\text{subject to}}} \text{tr } X + \text{tr } Y \leq 2 \quad (1858) \\ &\quad A \in \mathcal{C} \end{aligned}$$

- For  $A \in \mathbb{S}_+^N$  and  $\beta \in \mathbb{R}$

$$\begin{aligned} \beta \text{ tr } A &= \underset{\substack{X \in \mathbb{S}^N \\ \text{subject to}}} {\underset{\text{maximize}}{\text{maximize}}} \text{tr}(XA) \\ &\quad \text{subject to } X \preceq \beta I \end{aligned} \quad (1859)$$

But the following statement is numerically stable, preventing an unbounded solution in direction of a 0 eigenvalue:

$$\begin{aligned} \underset{\substack{X \in \mathbb{S}^N \\ \text{subject to}}} {\underset{\text{maximize}}{\text{maximize}}} &\quad \text{sgn}(\beta) \text{ tr}(XA) \\ &\quad \text{subject to } X \preceq |\beta| I \\ &\quad X \succeq -|\beta| I \end{aligned} \quad (1860)$$

where  $\beta \text{ tr } A = \text{tr}(X^* A)$ . If  $\beta \geq 0$ , then  $(X \succeq -|\beta| I) \leftarrow (X \succeq 0)$ .

---

<sup>C.2</sup> Hint: Write  $A = S\Sigma Q^T \in \mathbb{R}^{m \times n}$  and

$$\left[ \begin{array}{cc} X & A \\ A^T & Y \end{array} \right] = \left[ \begin{array}{c} U \\ V \end{array} \right] \left[ \begin{array}{cc} U^T & V^T \end{array} \right] \succeq 0$$

Show  $U^* = S\sqrt{\Sigma} \in \mathbb{R}^{m \times \min\{m, n\}}$  and  $V^* = Q\sqrt{\Sigma} \in \mathbb{R}^{n \times \min\{m, n\}}$ , hence  $\|U^*\|_F^2 = \|V^*\|_F^2$ .

- For symmetric  $A \in \mathbb{S}^N$ , its smallest and largest eigenvalue in  $\lambda(A) \in \mathbb{R}^N$  are respectively [12, §4.1] [44, §I.6.15] [228, §4.2] [268, §2.1] [269]

$$\min_i \{\lambda(A)_i\} = \inf_{\|x\|=1} x^T A x = \underset{\substack{x \in \mathbb{S}_+^N \\ \text{subject to } \text{tr } X = 1}}{\text{minimize}} \text{tr}(XA) = \underset{t \in \mathbb{R}}{\text{maximize}} t \quad (1861)$$

$$\max_i \{\lambda(A)_i\} = \sup_{\|x\|=1} x^T A x = \underset{\substack{x \in \mathbb{S}_+^N \\ \text{subject to } \text{tr } X = 1}}{\text{maximize}} \text{tr}(XA) = \underset{t \in \mathbb{R}}{\text{minimize}} t \quad (1862)$$

whereas

$$\lambda_N I \preceq A \preceq \lambda_1 I \quad (1863)$$

The largest eigenvalue  $\lambda_1$  is always convex in  $A \in \mathbb{S}^N$  because, given particular  $x$ ,  $x^T A x$  is linear in matrix  $A$ ; supremum of a family of linear functions is convex, as illustrated in Figure 78.C.3. So for  $A, B \in \mathbb{S}^N$ ,  $\lambda_1(A+B) \leq \lambda_1(A) + \lambda_1(B)$ . (1655) Similarly, the smallest eigenvalue  $\lambda_N$  of any symmetric matrix is a concave function of its entries;  $\lambda_N(A+B) \geq \lambda_N(A) + \lambda_N(B)$ . (1655) For  $v_N$  a normalized eigenvector of  $A$  corresponding to the smallest eigenvalue, and  $v_1$  a normalized eigenvector corresponding to the largest (principal) eigenvalue, the *principal eigenvector*,

$$v_N = \arg \inf_{\|x\|=1} x^T A x \quad (1864)$$

$$v_1 = \arg \sup_{\|x\|=1} x^T A x \quad (1865)$$

- For  $A \in \mathbb{S}^N$  having eigenvalues  $\lambda(A) \in \mathbb{R}^N$ , consider the unconstrained nonconvex optimization that is a projection of  $A$  on the rank-1 subset (§2.9.2.1, §3.6.0.0.1) of the boundary of positive semidefinite cone  $\mathbb{S}_+^N$ : Defining  $\lambda_1 \triangleq \max_i \{\lambda(A)_i\}$  and corresponding eigenvector  $v_1$

$$\begin{aligned} \underset{x}{\text{minimize}} \|xx^T - A\|_F^2 &= \underset{x}{\text{minimize}} \text{tr}(xx^T(x^Tx) - 2Axx^T + A^TA) \\ &= \begin{cases} \|\lambda(A)\|^2, & \lambda_1 \leq 0 \\ \|\lambda(A)\|^2 - \lambda_1^2, & \lambda_1 > 0 \end{cases} \end{aligned} \quad (1866)$$

$$\arg \underset{x}{\text{minimize}} \|xx^T - A\|_F^2 = \begin{cases} \mathbf{0}, & \lambda_1 \leq 0 \\ v_1 \sqrt{\lambda_1}, & \lambda_1 > 0 \end{cases} \quad (1867)$$

**Proof.** This is simply the Eckart & Young solution from §7.1.2:

$$x^* x^{*\top} = \begin{cases} \mathbf{0}, & \lambda_1 \leq 0 \\ \lambda_1 v_1 v_1^T, & \lambda_1 > 0 \end{cases} \quad (1868)$$

Given nonincreasingly ordered diagonalization  $A = Q\Lambda Q^T$  where  $\Lambda = \delta(\lambda(A))$  (§A.5), then (1866) has minimum value

$$\underset{x}{\text{minimize}} \|xx^T - A\|_F^2 = \begin{cases} \|Q\Lambda Q^T\|_F^2 = \|\delta(\Lambda)\|^2, & \lambda_1 \leq 0 \\ \left\| Q \begin{pmatrix} \lambda_1 & & \\ & 0 & \\ & & \ddots & \\ & & & 0 \end{pmatrix} - \Lambda \right\|_F^2 = \left\| \begin{pmatrix} \lambda_1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} - \delta(\Lambda) \right\|_F^2, & \lambda_1 > 0 \end{cases} \quad (1869)$$

---

<sup>C.3</sup>Largest eigenvalue  $\lambda_1$  is analogous to supremum over dashed vertical line segment in the figure. ♦

**C.2.0.0.2 Exercise.** *Rank-1 approximation.*

Given symmetric matrix  $A \in \mathbb{S}^N$ , prove:

$$\begin{aligned} v_N &= \arg \underset{x}{\text{maximize}} \quad \|xx^T - A\|_F^2 \\ &\text{subject to} \quad \|x\| = 1 \end{aligned} \quad (1870)$$

$$\begin{aligned} v_1 &= \arg \underset{x}{\text{minimize}} \quad \|xx^T - A\|_F^2 \\ &\text{subject to} \quad \|x\| = 1 \end{aligned} \quad (1871)$$

where  $v_N$  is a normalized eigenvector of  $A$  corresponding to its smallest eigenvalue and  $v_1$  corresponds to its largest. What is each objective's optimal value?  $\blacktriangledown$

- (Ky Fan, 1949) For eigenvalues  $\lambda(B) \in \mathbb{R}^N$  of  $B \in \mathbb{S}^N$  arranged in nonincreasing order, and for  $1 \leq k \leq N$  [12, §4.1] [242] [228, §4.3.18] [400, §2] [268, §2.1] [269]

$$\begin{aligned} \sum_{i=N-k+1}^N \lambda(B)_i &= \inf_{\substack{U \in \mathbb{R}^{N \times k} \\ U^T U = I}} \text{tr}(UU^T B) = \underset{\substack{X \in \mathbb{S}_+^N \\ \text{subject to}}}{} \text{minimize} \quad \text{tr}(XB) \\ &\quad X \preceq I \\ &\quad \text{tr } X = k \\ &= \underset{\substack{\mu \in \mathbb{R}, Z \in \mathbb{S}_+^N \\ \text{subject to}}}{} \text{maximize} \quad \mu(k-N) + \text{tr}(B-Z) \quad (b) \\ &\quad \mu I \succeq B - Z \quad (1872) \\ \sum_{i=1}^k \lambda(B)_i &= \sup_{\substack{U \in \mathbb{R}^{N \times k} \\ U^T U = I}} \text{tr}(UU^T B) = \underset{\substack{X \in \mathbb{S}_+^N \\ \text{subject to}}}{} \text{maximize} \quad \text{tr}(XB) \quad (c) \\ &\quad X \preceq I \\ &\quad \text{tr } X = k \\ &= \underset{\substack{\mu \in \mathbb{R}, Z \in \mathbb{S}_+^N \\ \text{subject to}}}{} \text{minimize} \quad \mu k + \text{tr } Z \quad (d) \\ &\quad \mu I \succeq B - Z \end{aligned}$$

Given ordered diagonalization  $B = Q\Lambda Q^T$ , (§A.5.1) then an optimal  $U$  for the infimum is  $U^* = Q(:, N-k+1:N) \in \mathbb{R}^{N \times k}$  whereas  $U^* = Q(:, 1:k) \in \mathbb{R}^{N \times k}$  for the supremum is more reliably computed. In both cases,  $X^* = U^*U^{*\top}$ . Optimization (a) searches the convex hull of outer product  $UU^T$  of all  $N \times k$  orthonormal matrices. (§2.3.2.0.1)

- For  $B \in \mathbb{S}^N$  whose eigenvalues  $\lambda(B) \in \mathbb{R}^N$  are arranged in nonincreasing order, and for diagonal matrix  $\Upsilon \in \mathbb{S}^k$  whose diagonal entries are arranged in nonincreasing order where  $1 \leq k \leq N$ , we utilize the main-diagonal  $\delta$  operator's selfadjointness property (1572): [13, §4.2]

$$\begin{aligned} \sum_{i=1}^k \Upsilon_{ii} \lambda(B)_{N-i+1} &= \inf_{\substack{U \in \mathbb{R}^{N \times k} \\ U^T U = I}} \text{tr}(\Upsilon U^T B U) = \inf_{\substack{U \in \mathbb{R}^{N \times k} \\ U^T U = I}} \delta(\Upsilon)^T \delta(U^T B U) \quad (1873) \\ &= \underset{\substack{V_i \in \mathbb{S}^N \\ \text{subject to}}}{} \text{minimize} \quad \text{tr} \left( B \sum_{i=1}^k (\Upsilon_{ii} - \Upsilon_{i+1,i+1}) V_i \right) \\ &\quad \text{tr } V_i = i, \quad i = 1 \dots k \\ &\quad I \succeq V_i \succeq 0, \quad i = 1 \dots k \end{aligned}$$

where  $\Upsilon_{k+1,k+1} \triangleq 0$ . We speculate,

$$\begin{aligned} \sum_{i=1}^k \Upsilon_{ii} \lambda(B)_i &= \sup_{\substack{U \in \mathbb{R}^{N \times k} \\ U^T U = I}} \text{tr}(\Upsilon U^T B U) = \sup_{\substack{U \in \mathbb{R}^{N \times k} \\ U^T U = I}} \delta(\Upsilon)^T \delta(U^T B U) \quad (1874) \end{aligned}$$

Alizadeh shows: [12, §4.2]

$$\begin{aligned}
 \sum_{i=1}^k \Upsilon_{ii} \lambda(B)_i &= \underset{\mu \in \mathbb{R}^k, Z_i \in \mathbb{S}^N}{\text{minimize}} \quad \sum_{i=1}^k i \mu_i + \text{tr } Z_i \\
 &\text{subject to} \quad \begin{aligned} \mu_i I + Z_i - (\Upsilon_{ii} - \Upsilon_{i+1,i+1})B &\succeq 0, & i = 1 \dots k \\ Z_i \succeq 0, & & i = 1 \dots k \end{aligned} \\
 &= \underset{V_i \in \mathbb{S}^N}{\text{maximize}} \quad \text{tr} \left( B \sum_{i=1}^k (\Upsilon_{ii} - \Upsilon_{i+1,i+1}) V_i \right) \\
 &\text{subject to} \quad \begin{aligned} \text{tr } V_i = i, & & i = 1 \dots k \\ I \succeq V_i \succeq 0, & & i = 1 \dots k \end{aligned} \tag{1875}
 \end{aligned}$$

where  $\Upsilon_{k+1,k+1} \triangleq 0$ .

- The largest eigenvalue magnitude  $\mu$  of  $A \in \mathbb{S}^N$

$$\begin{aligned}
 \max_i \{|\lambda(A)_i|\} &= \underset{\mu \in \mathbb{R}}{\text{minimize}} \quad \mu \\
 &\text{subject to} \quad -\mu I \preceq A \preceq \mu I \tag{1876}
 \end{aligned}$$

is minimized over convex set  $\mathcal{C}$  by semidefinite program: (confer §7.1.5)

$$\begin{aligned}
 \underset{\substack{A \\ \text{subject to} \\ A \in \mathcal{C}}}{\text{minimize}} \quad \|A\|_2 &\equiv \underset{\substack{\mu, A \\ \text{subject to} \\ -\mu I \preceq A \preceq \mu I \\ A \in \mathcal{C}}}{\text{minimize}} \quad \mu \tag{1877}
 \end{aligned}$$

id est,

$$\mu^* \triangleq \max_i \{|\lambda(A^*)_i|, i = 1 \dots N\} \in \mathbb{R}_+ \tag{1878}$$

- For  $B \in \mathbb{S}^N$  whose eigenvalues  $\lambda(B) \in \mathbb{R}^N$  are arranged in nonincreasing order, let  $\Pi \lambda(B)$  be a permutation of eigenvalues  $\lambda(B)$  such that their absolute value becomes arranged in nonincreasing order:  $|\Pi \lambda(B)|_1 \geq |\Pi \lambda(B)|_2 \geq \dots \geq |\Pi \lambda(B)|_N$ . Then, for  $1 \leq k \leq N$  [12, §4.3] C.4

$$\begin{aligned}
 \sum_{i=1}^k |\Pi \lambda(B)|_i &= \underset{\mu \in \mathbb{R}, Z \in \mathbb{S}_+^N}{\text{minimize}} \quad k \mu + \text{tr } Z = \underset{\substack{V, W \in \mathbb{S}_+^N \\ \text{subject to} \\ I \succeq V, W \\ \text{tr}(V+W)=k}}{\text{maximize}} \quad \langle B, V - W \rangle \\
 &\text{subject to} \quad \begin{aligned} \mu I + Z + B &\succeq 0 \\ \mu I + Z - B &\succeq 0 \end{aligned} \tag{1879}
 \end{aligned}$$

For diagonal matrix  $\Upsilon \in \mathbb{S}^k$  whose diagonal entries are arranged in nonincreasing order where  $1 \leq k \leq N$

$$\begin{aligned}
 \sum_{i=1}^k \Upsilon_{ii} |\Pi \lambda(B)|_i &= \underset{\mu \in \mathbb{R}^k, Z_i \in \mathbb{S}^N}{\text{minimize}} \quad \sum_{i=1}^k i \mu_i + \text{tr } Z_i \\
 &\text{subject to} \quad \begin{aligned} \mu_i I + Z_i + (\Upsilon_{ii} - \Upsilon_{i+1,i+1})B &\succeq 0, & i = 1 \dots k \\ \mu_i I + Z_i - (\Upsilon_{ii} - \Upsilon_{i+1,i+1})B &\succeq 0, & i = 1 \dots k \\ Z_i \succeq 0, & & i = 1 \dots k \end{aligned} \\
 &= \underset{\substack{V_i, W_i \in \mathbb{S}^N \\ \text{subject to} \\ \text{tr}(V_i + W_i) = i \\ I \succeq V_i \succeq 0 \\ I \succeq W_i \succeq 0}}{\text{maximize}} \quad \text{tr} \left( B \sum_{i=1}^k (\Upsilon_{ii} - \Upsilon_{i+1,i+1})(V_i - W_i) \right) \\
 &\text{subject to} \quad \begin{aligned} \text{tr } V_i = i, & & i = 1 \dots k \\ I \succeq V_i \succeq 0, & & i = 1 \dots k \\ I \succeq W_i \succeq 0, & & i = 1 \dots k \end{aligned} \tag{1880}
 \end{aligned}$$

where  $\Upsilon_{k+1,k+1} \triangleq 0$ .

---

C.4 We eliminate a redundant positive semidefinite variable from Alizadeh's minimization. There exist typographical errors in [326, (6.49) (6.55)] for this minimization.

**C.2.0.0.3 Exercise.** *Weighted sum of largest eigenvalues.*

Prove (1874). ▼

- For  $A, B \in \mathbb{S}^N$  whose eigenvalues  $\lambda(A), \lambda(B) \in \mathbb{R}^N$  are respectively arranged in nonincreasing order, and for nonincreasingly ordered diagonalizations  $A = W_A \Upsilon W_A^T$  and  $B = W_B \Lambda W_B^T$  [226] [268, §2.1] [269]

$$\lambda(A)^T \lambda(B) = \sup_{\substack{U \in \mathbb{R}^{N \times N} \\ U^T U = I}} \text{tr}(A^T U^T B U) \geq \text{tr}(A^T B) \quad (1899)$$

(confer (1904)) where optimal  $U$  is

$$U^* = W_B W_A^T \in \mathbb{R}^{N \times N} \quad (1896)$$

We can push that upper bound higher using a result in §C.4.2.1:

$$|\lambda(A)|^T |\lambda(B)| = \sup_{\substack{U \in \mathbb{C}^{N \times N} \\ U^H U = I}} \text{re tr}(A^T U^H B U) \quad (1881)$$

For step function  $\psi$  as defined in (1733), optimal  $U$  becomes

$$U^* = W_B \sqrt{\delta(\psi(\delta(\Lambda)))}^H \sqrt{\delta(\psi(\delta(\Upsilon)))} W_A^T \in \mathbb{C}^{N \times N} \quad (1882)$$

### C.3 Orthogonal Procrustes problem

Given matrices  $A, B \in \mathbb{R}^{n \times N}$ , their product having full singular value decomposition (§A.6.3)

$$AB^T \triangleq U \Sigma Q^T \in \mathbb{R}^{n \times n} \quad (1883)$$

then an optimal solution  $R^*$  to the orthogonal Procrustes problem

$$\begin{aligned} & \underset{R}{\text{minimize}} \quad \|A - R^T B\|_F \\ & \text{subject to} \quad R^T = R^{-1} \end{aligned} \quad (1884)$$

maximizes  $\text{tr}(A^T R^T B)$  over the nonconvex manifold of orthogonal matrices: [228, §7.4.8]

$$R^* = Q U^T \in \mathbb{R}^{n \times n} \quad (1885)$$

A necessary and sufficient condition for optimality

$$AB^T R^* \succeq 0 \quad (1886)$$

holds whenever  $R^*$  is an orthogonal matrix. [188, §4]Optimal solution  $R^*$  can reveal rotation/reflection (§5.5.2, §B.5) of one list in the columns of matrix  $A$  with respect to another list in  $B$ . Solution is unique if  $\text{rank } BV_N = n$ . [126, §2.4.1] In the case that  $A$  is a vector and permutation of  $B$ , solution  $R^*$  is not necessarily a permutation matrix (§4.7.0.0.3) although the optimal objective will be zero. More generally, the optimal value for objective of minimization is

$$\begin{aligned} \text{tr}(A^T A + B^T B - 2AB^T R^*) &= \text{tr}(A^T A) + \text{tr}(B^T B) - 2 \text{tr}(U \Sigma U^T) \\ &= \|A\|_F^2 + \|B\|_F^2 - 2\delta(\Sigma)^T \mathbf{1} \end{aligned} \quad (1887)$$

while the optimal value for corresponding trace maximization is

$$\sup_{R^T = R^{-1}} \text{tr}(A^T R^T B) = \text{tr}(A^T R^{*T} B) = \delta(\Sigma)^T \mathbf{1} \geq \text{tr}(A^T B) \quad (1888)$$

The same optimal solution  $R^*$  solves

$$\begin{aligned} & \underset{R}{\text{maximize}} \quad \|A + R^T B\|_F \\ & \text{subject to} \quad R^T = R^{-1} \end{aligned} \quad (1889)$$

### C.3.0.1 Procrustes relaxation

By replacing its feasible set with (Example 2.3.2.0.5) the convex hull of orthogonal matrices, we relax Procrustes problem (1884) to a convex problem

$$\begin{array}{ll} \text{minimize}_R & \|A - R^T B\|_F^2 = \text{tr}(A^T A + B^T B) - 2 \text{maximize}_R \text{tr}(A^T R^T B) \\ \text{subject to} & R^T = R^{-1} \quad \text{subject to} \quad \begin{bmatrix} I & R \\ R^T & I \end{bmatrix} \succeq 0 \end{array} \quad (1890)$$

whose adjusted objective must always equal Procrustes'.<sup>C.5</sup>

### C.3.1 Effect of translation

Consider the impact on problem (1884) of DC offset in known lists  $A, B \in \mathbb{R}^{n \times N}$ . Rotation of  $B$  there is with respect to the origin, so better results may be obtained if offset is first accounted. Because geometric centers of lists  $AV$  and  $BV$  are the origin, instead we solve

$$\begin{array}{ll} \text{minimize}_R & \|AV - R^T BV\|_F \\ \text{subject to} & R^T = R^{-1} \end{array} \quad (1891)$$

where  $V \in \mathbb{S}^N$  is the geometric centering matrix (§B.4.1). Now we define the full singular value decomposition

$$AVB^T \triangleq U\Sigma Q^T \in \mathbb{R}^{n \times n} \quad (1892)$$

and an optimal rotation matrix

$$R^* = QU^T \in \mathbb{R}^{n \times n} \quad (1885)$$

The desired result is an optimally rotated offset list

$$R^{*T}BV + A(I - V) \approx A \quad (1893)$$

which most closely matches the list in  $A$ . Equality is attained when the lists are precisely related by a rotation/reflection and an offset. When  $R^{*T}B = A$  or  $B\mathbf{1} = A\mathbf{1} = \mathbf{0}$ , this result (1893) reduces to  $R^{*T}B \approx A$ .

#### C.3.1.1 Translation of extended list

Suppose an optimal rotation matrix  $R^* \in \mathbb{R}^{n \times n}$  were derived as before from matrix  $B \in \mathbb{R}^{n \times N}$ , but  $B$  is part of a larger list in the columns of  $[C \ B] \in \mathbb{R}^{n \times M+N}$  where  $C \in \mathbb{R}^{n \times M}$ . In that event, we wish to apply the rotation/reflection and translation to the larger list. The expression supplanting the approximation in (1893) makes  $\mathbf{1}^T$  of compatible dimension;

$$R^{*T}[C - B\mathbf{1}\mathbf{1}^T \frac{1}{N} \quad BV] + A\mathbf{1}\mathbf{1}^T \frac{1}{N} \quad (1894)$$

*id est*,  $C - B\mathbf{1}\mathbf{1}^T \frac{1}{N} \in \mathbb{R}^{n \times M}$  and  $A\mathbf{1}\mathbf{1}^T \frac{1}{N} \in \mathbb{R}^{n \times M+N}$ .

---

<sup>C.5</sup> (because orthogonal matrices are the extreme points of this hull) and whose optimal numerical solution (SDPT3 [389]) [191] is reliably observed to be orthogonal for  $n \leq N$ .

## C.4 Two-sided orthogonal Procrustes

### C.4.0.1 Minimization

Given symmetric  $A, B \in \mathbb{S}^N$ , each having diagonalization (§A.5.1)

$$A \triangleq Q_A \Lambda_A Q_A^T, \quad B \triangleq Q_B \Lambda_B Q_B^T \quad (1895)$$

where eigenvalues are arranged in their respective diagonal matrix  $\Lambda$  in nonincreasing order, then an optimal solution [148]

$$R^* = Q_B Q_A^T \in \mathbb{R}^{N \times N} \quad (1896)$$

to the two-sided orthogonal Procrustes problem

$$\begin{array}{ll} \underset{R}{\text{minimize}} & \|A - R^T B R\|_F \\ \text{subject to} & R^T = R^{-1} \end{array} = \begin{array}{ll} \underset{R}{\text{minimize}} & \text{tr}(A^T A - 2A^T R^T B R + B^T B) \\ \text{subject to} & R^T = R^{-1} \end{array} \quad (1897)$$

maximizes  $\text{tr}(A^T R^T B R)$  over the nonconvex manifold of orthogonal matrices. Optimal product  $R^{*T} B R^*$  has the eigenvectors of  $A$  but the eigenvalues of  $B$ . [188, §7.5.1] The optimal value for the objective of minimization is, by (49)

$$\|Q_A \Lambda_A Q_A^T - R^{*T} Q_B \Lambda_B Q_B^T R^*\|_F = \|Q_A (\Lambda_A - \Lambda_B) Q_A^T\|_F = \|\Lambda_A - \Lambda_B\|_F \quad (1898)$$

while the corresponding trace maximization has optimal value

$$\sup_{R^T=R^{-1}} \text{tr}(A^T R^T B R) = \text{tr}(A^T R^{*T} B R^*) = \text{tr}(\Lambda_A \Lambda_B) \geq \text{tr}(A^T B) \quad (1899)$$

The lower bound on inner product of eigenvalues is due to Fan (p.493).

### C.4.0.2 Maximization

Any permutation matrix is an orthogonal matrix. Defining a row- and column-swapping permutation matrix (a reflection matrix, §B.5.3)

$$\Xi = \Xi^T = \begin{bmatrix} \mathbf{0} & & & 1 \\ & \ddots & & \\ & & 1 & \\ 1 & & & \mathbf{0} \end{bmatrix} \quad (1900)$$

then an optimal solution  $R^*$  to the maximization problem [*sic*]

$$\begin{array}{ll} \underset{R}{\text{maximize}} & \|A - R^T B R\|_F \\ \text{subject to} & R^T = R^{-1} \end{array} \quad (1901)$$

minimizes  $\text{tr}(A^T R^T B R)$ : [226] [268, §2.1] [269]

$$R^* = Q_B \Xi Q_A^T \in \mathbb{R}^{N \times N} \quad (1902)$$

The optimal value for the objective of maximization is

$$\begin{aligned} \|Q_A \Lambda_A Q_A^T - R^{*T} Q_B \Lambda_B Q_B^T R^*\|_F &= \|Q_A \Lambda_A Q_A^T - Q_A \Xi^T \Lambda_B \Xi Q_A^T\|_F \\ &= \|\Lambda_A - \Xi \Lambda_B \Xi\|_F \end{aligned} \quad (1903)$$

while the corresponding trace minimization has optimal value

$$\inf_{R^T=R^{-1}} \text{tr}(A^T R^T B R) = \text{tr}(A^T R^{*T} B R^*) = \text{tr}(\Lambda_A \Xi \Lambda_B \Xi) \quad (1904)$$

### C.4.1 Procrustes' relation to linear programming

Although these two-sided Procrustes problems are nonconvex, there is a connection with linear programming [13, §3] [268, §2.1] [269]: Given  $A, B \in \mathbb{S}^N$ , this semidefinite program in  $S$  and  $T$

$$\begin{array}{ll} \underset{R}{\text{minimize}} & \text{tr}(A^T R^T B R) = \underset{S, T \in \mathbb{S}^N}{\text{maximize}} \text{tr}(S + T) \\ \text{subject to} & R^T = R^{-1} \quad \text{subject to} \quad A^T \otimes B - I \otimes S - T \otimes I \succeq 0 \end{array} \quad (1905)$$

(where  $\otimes$  signifies Kronecker product (§D.1.2.1)) has optimal objective value (1904). These two problems in (1905) are strong duals (§2.13.1.1.1). Given ordered diagonalizations (1895), make the observation:

$$\inf_R \text{tr}(A^T R^T B R) = \inf_{\hat{R}} \text{tr}(\Lambda_A \hat{R}^T \Lambda_B \hat{R}) \quad (1906)$$

because  $\hat{R} \triangleq Q_B^T R Q_A$  on the set of orthogonal matrices (which includes the permutation matrices) is a bijection. This means, basically, diagonal matrices of eigenvalues  $\Lambda_A$  and  $\Lambda_B$  may be substituted for  $A$  and  $B$ , so only the main diagonals of  $S$  and  $T$  come into play;

$$\begin{array}{ll} \underset{S, T \in \mathbb{S}^N}{\text{maximize}} & \mathbf{1}^T \delta(S + T) \\ \text{subject to} & \delta(\Lambda_A \otimes (\Xi \Lambda_B \Xi) - I \otimes S - T \otimes I) \succeq 0 \end{array} \quad (1907)$$

a linear program in  $\delta(S)$  and  $\delta(T)$  having the same optimal objective value as the semidefinite program (1905).

We relate their results to Procrustes problem (1897) by manipulating signs (1843) and permuting eigenvalues:

$$\begin{array}{ll} \underset{R}{\text{maximize}} & \text{tr}(A^T R^T B R) = \underset{S, T \in \mathbb{S}^N}{\text{minimize}} \mathbf{1}^T \delta(S + T) \\ \text{subject to} & R^T = R^{-1} \quad \text{subject to} \quad \delta(I \otimes S + T \otimes I - \Lambda_A \otimes \Lambda_B) \succeq 0 \\ & = \underset{S, T \in \mathbb{S}^N}{\text{minimize}} \text{tr}(S + T) \\ & \text{subject to} \quad I \otimes S + T \otimes I - A^T \otimes B \succeq 0 \end{array} \quad (1908)$$

This formulation has optimal objective value identical to that in (1899).

### C.4.2 Two-sided orthogonal Procrustes via SVD

By making left- and right-side orthogonal matrices independent, we can push the upper bound on trace (1899) a little further: Given real matrices  $A, B$  each having full singular value decomposition (§A.6.3)

$$A \triangleq U_A \Sigma_A Q_A^T \in \mathbb{R}^{m \times n}, \quad B \triangleq U_B \Sigma_B Q_B^T \in \mathbb{R}^{m \times n} \quad (1909)$$

then a well-known optimal solution  $R^*, S^*$  to the problem

$$\begin{array}{ll} \underset{R, S}{\text{minimize}} & \|A - SBR\|_F \\ \text{subject to} & R^H = R^{-1} \\ & S^H = S^{-1} \end{array} \quad (1910)$$

maximizes  $\text{retr}(A^T S B R)$ : [351] [320] [55] [220] optimal orthogonal matrices

$$S^* = U_A U_B^H \in \mathbb{R}^{m \times m}, \quad R^* = Q_B Q_A^H \in \mathbb{R}^{n \times n} \quad (1911)$$

[sic] are not necessarily unique [228, §7.4.13] because the feasible set is not convex. The optimal value for the objective of minimization is, by (49)

$$\|U_A \Sigma_A Q_A^H - S^* U_B \Sigma_B Q_B^H R^*\|_F = \|U_A (\Sigma_A - \Sigma_B) Q_A^H\|_F = \|\Sigma_A - \Sigma_B\|_F \quad (1912)$$

while the corresponding trace maximization has optimal value [44, §III.6.12]

$$\sup_{\substack{R^H = R^{-1} \\ S^H = S^{-1}}} |\text{tr}(A^T S B R)| = \sup_{\substack{R^H = R^{-1} \\ S^H = S^{-1}}} \text{re tr}(A^T S B R) = \text{re tr}(A^T S^* B R^*) = \text{tr}(\Sigma_A^T \Sigma_B) \geq \text{tr}(A^T B) \quad (1913)$$

for which it is necessary

$$A^T S^* B R^* \succeq 0, \quad B R^* A^T S^* \succeq 0 \quad (1914)$$

The lower bound on inner product of singular values in (1913) is due to von Neumann. Equality is attained if  $U_A^H U_B = I$  and  $Q_B^H Q_A = I$ .

#### C.4.2.1 Symmetric matrices

Now optimizing over the complex manifold of unitary matrices (§B.5.2), the upper bound on trace (1899) is thereby raised: Suppose we are given diagonalizations for (real) symmetric  $A, B$  (§A.5)

$$A = W_A \Upsilon W_A^T \in \mathbb{S}^n, \quad \delta(\Upsilon) \in \mathcal{K}_M \quad (1915)$$

$$B = W_B \Lambda W_B^T \in \mathbb{S}^n, \quad \delta(\Lambda) \in \mathcal{K}_M \quad (1916)$$

having their respective eigenvalues in diagonal matrices  $\Upsilon, \Lambda \in \mathbb{S}^n$  arranged in nonincreasing order (membership to the monotone cone  $\mathcal{K}_M$  (441)). Then by splitting eigenvalue signs, we invent a symmetric SVD-like decomposition

$$A \triangleq U_A \Sigma_A Q_A^H \in \mathbb{S}^n, \quad B \triangleq U_B \Sigma_B Q_B^H \in \mathbb{S}^n \quad (1917)$$

where  $U_A, U_B, Q_A, Q_B \in \mathbb{C}^{n \times n}$  are unitary matrices defined by (confer §A.6.3.2)

$$U_A \triangleq W_A \sqrt{\delta(\psi(\delta(\Upsilon)))}, \quad Q_A \triangleq W_A \sqrt{\delta(\psi(\delta(\Upsilon)))}^H, \quad \Sigma_A = |\Upsilon| \quad (1918)$$

$$U_B \triangleq W_B \sqrt{\delta(\psi(\delta(\Lambda)))}, \quad Q_B \triangleq W_B \sqrt{\delta(\psi(\delta(\Lambda)))}^H, \quad \Sigma_B = |\Lambda| \quad (1919)$$

where step function  $\psi$  is defined in (1733). In this circumstance,

$$S^* = U_A U_B^H = R^{*T} \in \mathbb{C}^{n \times n} \quad (1920)$$

optimal matrices (1911) now unitary are related by transposition. The optimal value of objective (1912) is

$$\|U_A \Sigma_A Q_A^H - S^* U_B \Sigma_B Q_B^H R^*\|_F = \|\Upsilon - \Lambda\|_F \quad (1921)$$

while the corresponding optimal value of trace maximization (1913) is

$$\sup_{\substack{R^H = R^{-1} \\ S^H = S^{-1}}} \text{re tr}(A^T S B R) = \text{tr}(|\Upsilon| |\Lambda|) \quad (1922)$$

### C.4.2.2 Diagonal matrices

Now suppose  $A$  and  $B$  are diagonal matrices

$$A = \Upsilon = \delta^2(\Upsilon) \in \mathbb{S}^n, \quad \delta(\Upsilon) \in \mathcal{K}_{\mathcal{M}} \quad (1923)$$

$$B = \Lambda = \delta^2(\Lambda) \in \mathbb{S}^n, \quad \delta(\Lambda) \in \mathcal{K}_{\mathcal{M}} \quad (1924)$$

both having their respective main diagonal entries arranged in nonincreasing order:

$$\begin{aligned} & \underset{R, S}{\text{minimize}} \quad \|\Upsilon - S\Lambda R\|_{\text{F}} \\ & \text{subject to} \quad R^H = R^{-1} \\ & \qquad \qquad \qquad S^H = S^{-1} \end{aligned} \quad (1925)$$

Then we have a symmetric decomposition from unitary matrices as in (1917) where

$$U_A \triangleq \sqrt{\delta(\psi(\delta(\Upsilon)))}, \quad Q_A \triangleq \sqrt{\delta(\psi(\delta(\Upsilon)))}^H, \quad \Sigma_A = |\Upsilon| \quad (1926)$$

$$U_B \triangleq \sqrt{\delta(\psi(\delta(\Lambda)))}, \quad Q_B \triangleq \sqrt{\delta(\psi(\delta(\Lambda)))}^H, \quad \Sigma_B = |\Lambda| \quad (1927)$$

Procrustes solution (1911) again sees the transposition relationship

$$S^* = U_A U_B^H = R^{*\top} \in \mathbb{C}^{n \times n} \quad (1920)$$

but both optimal unitary matrices are now themselves diagonal. So,

$$S^* \Lambda R^* = \delta(\psi(\delta(\Upsilon))) \Lambda \delta(\psi(\delta(\Lambda))) = \delta(\psi(\delta(\Upsilon))) |\Lambda| \quad (1928)$$

## C.5 Quadratics

### C.5.1 minimization, convex

Given positive semidefinite matrix  $A \succeq 0$  (§A.4.0.0.2)

$$\inf_{x \in \mathbb{R}^n} \frac{1}{2} x^T A x + b^T x = \frac{1}{2} \inf_{x \in \mathbb{R}^n} [x^T \ 1] \begin{bmatrix} A & b \\ b^T & 0 \end{bmatrix} \begin{bmatrix} x \\ 1 \end{bmatrix} = \begin{cases} -\frac{1}{2} b^T A^\dagger b, & b \in \mathcal{R}(A) \\ -\infty, & \text{otherwise} \end{cases} \quad (1929)$$

where  $b \in \mathcal{R}(A)$  is condition (1664) of the Schur complement.

#### C.5.1.0.1 Exercise. maximization, convex case.

Assume a negative semidefinite matrix  $A \preceq 0$ . Write the analogue to (1929) for supremum of a concave quadratic. ▼

### C.5.2 minimization, nonconvex

[385, §2] [358, §2] Given symmetric matrix  $A \in \mathbb{S}^n$ , vector  $b \in \mathbb{R}^n$ , and scalar  $\rho > 0$

$$\begin{aligned} & \underset{x}{\text{minimize}} \quad \frac{1}{2} x^T A x + b^T x \quad \Leftrightarrow \quad \begin{aligned} \text{i)} \quad & (A + \lambda^* I)x^* = -b \\ \text{subject to} \quad & \|x\| \leq \rho \quad \text{ii)} \quad \lambda^*(\|x^*\| - \rho) = 0, \quad \|x^*\| \leq \rho \\ & \text{iii)} \quad A + \lambda^* I \succeq 0 \end{aligned} \end{aligned} \quad (1930)$$

is a nonconvex problem for symmetric  $A$  unless  $A \succeq 0$ . But necessary and sufficient global optimality conditions are known for any symmetric  $A$ : vector  $x^*$  solves minimization (1930) iff  $\exists$  Lagrange multiplier  $\lambda^* \geq 0$  satisfying the three corresponding conditions.

Conditions **i** and **ii** are necessary KKT conditions, [65, §5.5.3] while condition **iii** governs passage to nonconvex global optimality and derived from (1929) like so: Lagrangian

$$\mathcal{L}(x, \lambda) = \frac{1}{2}x^T Ax + b^T x + \lambda(x^T x - \rho^2) = \frac{1}{2}x^T(A + 2\lambda I)x + b^T x - \lambda\rho^2 \quad (1931)$$

has finite infimum, assuming  $A + 2\lambda I \succeq 0$

$$\inf_{x \in \mathbb{R}^n} \mathcal{L}(x, \lambda) = -\frac{1}{2}b^T(A + 2\lambda I)^\dagger b - \rho^2\lambda, \quad b \in \mathcal{R}(A + 2\lambda I) \quad (1932)$$

that is a lower bound to generally nonconvex problem (1930).  $\lambda^*$  is unique; it is the solution to a convex dual problem that attempts the greatest lower bound to (1930), substituting  $\lambda \leftarrow \frac{1}{2}\lambda$

$$\begin{aligned} \underset{\lambda \in \mathbb{R}_+}{\text{maximize}} \quad & \inf_{x \in \mathbb{R}^n} \mathcal{L}(x, \lambda) = \underset{\lambda \in \mathbb{R}_+}{\text{maximize}} \quad -b^T(A + \lambda I)^\dagger b - \rho^2\lambda \\ & \text{subject to} \quad A + \lambda I \succeq 0 \\ & \quad b \in \mathcal{R}(A + \lambda I) \end{aligned} \quad (1933)$$

$x^*$  is unique if  $A + \lambda^* I \succ 0$ .

Equality-constrained problem

$$\begin{aligned} \underset{x}{\text{minimize}} \quad & \frac{1}{2}x^T Ax + b^T x \\ \text{subject to} \quad & \|x\| = \rho \end{aligned} \quad \Leftrightarrow \quad \begin{array}{ll} \text{i)} & (A + \lambda^* I)x^* = -b \\ \text{ii)} & \|x^*\| = \rho \\ \text{iii)} & A + \lambda^* I \succeq 0 \end{array} \quad (1934)$$

is nonconvex for any symmetric  $A$  matrix.  $x^*$  solves minimization (1934) iff  $\exists \lambda^* \in \mathbb{R}$  satisfying the associated conditions.  $\lambda^*$  and  $x^*$  are unique as before.

### C.5.3 maximization, nonconvex

Hiriart-Urruty disclosed global optimality conditions in 1998 [223]<sup>C.6</sup> for maximizing a convex quadratic with convex constraints; a nonconvex problem [343, §32].

---

<sup>C.6</sup>... the assumptions in Theorem 8 ask for the  $Q_i$  being positive definite (see the top of the page of Theorem 8). I must confess that I do not remember why. — Jean-Baptiste Hiriart-Urruty



# Appendix D

## Matrix calculus

*From too much study, and from extreme passion, cometh madnesse.*

— Isaac Newton [175, §5]

### D.1 Gradient, Directional derivative, Taylor series

#### D.1.1 Gradients

Gradient of a differentiable real function  $f(x) : \mathbb{R}^K \rightarrow \mathbb{R}$  with respect to its vector argument is defined uniquely in terms of partial derivatives

$$\nabla f(x) \triangleq \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_K} \end{bmatrix} \in \mathbb{R}^K \quad (1935)$$

while the second-order gradient of the twice differentiable real function with respect to its vector argument is traditionally called the *Hessian*;

$$\nabla^2 f(x) \triangleq \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_K} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_K \partial x_1} & \frac{\partial^2 f(x)}{\partial x_K \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_K^2} \end{bmatrix} \in \mathbb{S}^K \quad (1936)$$

The gradient of vector-valued function  $v(x) : \mathbb{R} \rightarrow \mathbb{R}^N$  on real domain is a row vector

$$\nabla v(x) \triangleq \left[ \frac{\partial v_1(x)}{\partial x} \quad \frac{\partial v_2(x)}{\partial x} \quad \cdots \quad \frac{\partial v_N(x)}{\partial x} \right] \in \mathbb{R}^N \quad (1937)$$

while the second-order gradient is

$$\nabla^2 v(x) \triangleq \left[ \frac{\partial^2 v_1(x)}{\partial x^2} \quad \frac{\partial^2 v_2(x)}{\partial x^2} \quad \cdots \quad \frac{\partial^2 v_N(x)}{\partial x^2} \right] \in \mathbb{R}^N \quad (1938)$$

Gradient of vector-valued function  $h(x) : \mathbb{R}^K \rightarrow \mathbb{R}^N$  on vector domain is

$$\nabla h(x) \triangleq \begin{bmatrix} \frac{\partial h_1(x)}{\partial x_1} & \frac{\partial h_2(x)}{\partial x_1} & \dots & \frac{\partial h_N(x)}{\partial x_1} \\ \frac{\partial h_1(x)}{\partial x_2} & \frac{\partial h_2(x)}{\partial x_2} & \dots & \frac{\partial h_N(x)}{\partial x_2} \\ \vdots & \vdots & & \vdots \\ \frac{\partial h_1(x)}{\partial x_K} & \frac{\partial h_2(x)}{\partial x_K} & \dots & \frac{\partial h_N(x)}{\partial x_K} \end{bmatrix} = [\nabla h_1(x) \ \nabla h_2(x) \ \dots \ \nabla h_N(x)] \in \mathbb{R}^{K \times N} \quad (1939)$$

while the second-order gradient has a three-dimensional written representation dubbed *cubix*<sup>D.1</sup>,

$$\nabla^2 h(x) \triangleq \begin{bmatrix} \nabla \frac{\partial h_1(x)}{\partial x_1} & \nabla \frac{\partial h_2(x)}{\partial x_1} & \dots & \nabla \frac{\partial h_N(x)}{\partial x_1} \\ \nabla \frac{\partial h_1(x)}{\partial x_2} & \nabla \frac{\partial h_2(x)}{\partial x_2} & \dots & \nabla \frac{\partial h_N(x)}{\partial x_2} \\ \vdots & \vdots & & \vdots \\ \nabla \frac{\partial h_1(x)}{\partial x_K} & \nabla \frac{\partial h_2(x)}{\partial x_K} & \dots & \nabla \frac{\partial h_N(x)}{\partial x_K} \end{bmatrix} = [\nabla^2 h_1(x) \ \nabla^2 h_2(x) \ \dots \ \nabla^2 h_N(x)] \in \mathbb{R}^{K \times N \times K} \quad (1940)$$

where the gradient of each real entry is with respect to vector  $x$  as in (1935).

The gradient of real function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$  on matrix domain is

$$\nabla g(X) \triangleq \begin{bmatrix} \frac{\partial g(X)}{\partial X_{11}} & \frac{\partial g(X)}{\partial X_{12}} & \dots & \frac{\partial g(X)}{\partial X_{1L}} \\ \frac{\partial g(X)}{\partial X_{21}} & \frac{\partial g(X)}{\partial X_{22}} & \dots & \frac{\partial g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial g(X)}{\partial X_{K1}} & \frac{\partial g(X)}{\partial X_{K2}} & \dots & \frac{\partial g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L} \quad (1941)$$

$$= \begin{bmatrix} \nabla_{X(:,1)} g(X) \\ \nabla_{X(:,2)} g(X) \\ \vdots \\ \nabla_{X(:,L)} g(X) \end{bmatrix} \in \mathbb{R}^{K \times 1 \times L}$$

where gradient  $\nabla_{X(:,i)}$  is with respect to the  $i^{\text{th}}$  column of  $X$ . The strange appearance of (1941) in  $\mathbb{R}^{K \times 1 \times L}$  is meant to suggest a third dimension perpendicular to the page (not a diagonal matrix). The second-order gradient has representation

---

<sup>D.1</sup>The word *matrix* comes from the Latin for *womb*; related to the prefix *matri-* derived from *mater* meaning *mother*.

$$\nabla^2 g(X) \triangleq \begin{bmatrix} \nabla \frac{\partial g(X)}{\partial X_{11}} & \nabla \frac{\partial g(X)}{\partial X_{12}} & \cdots & \nabla \frac{\partial g(X)}{\partial X_{1L}} \\ \nabla \frac{\partial g(X)}{\partial X_{21}} & \nabla \frac{\partial g(X)}{\partial X_{22}} & \cdots & \nabla \frac{\partial g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \nabla \frac{\partial g(X)}{\partial X_{K1}} & \nabla \frac{\partial g(X)}{\partial X_{K2}} & \cdots & \nabla \frac{\partial g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times K \times L} \quad (1942)$$

$$= \begin{bmatrix} \nabla \nabla_{X(:,1)} g(X) \\ \nabla \nabla_{X(:,2)} g(X) \\ \ddots \\ \nabla \nabla_{X(:,L)} g(X) \end{bmatrix} \in \mathbb{R}^{K \times 1 \times L \times K \times L}$$

where the gradient  $\nabla$  is with respect to matrix  $X$ .

Gradient of vector-valued function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^N$  on matrix domain is a cubix

$$\nabla g(X) \triangleq \begin{bmatrix} \nabla_{X(:,1)} g_1(X) & \nabla_{X(:,1)} g_2(X) & \cdots & \nabla_{X(:,1)} g_N(X) \\ \nabla_{X(:,2)} g_1(X) & \nabla_{X(:,2)} g_2(X) & \cdots & \nabla_{X(:,2)} g_N(X) \\ \ddots & \ddots & \ddots & \ddots \\ \nabla_{X(:,L)} g_1(X) & \nabla_{X(:,L)} g_2(X) & \cdots & \nabla_{X(:,L)} g_N(X) \end{bmatrix} \quad (1943)$$

$$= [\nabla g_1(X) \ \nabla g_2(X) \ \cdots \ \nabla g_N(X)] \in \mathbb{R}^{K \times N \times L}$$

while the second-order gradient has a five-dimensional representation;

$$\nabla^2 g(X) \triangleq \begin{bmatrix} \nabla \nabla_{X(:,1)} g_1(X) & \nabla \nabla_{X(:,1)} g_2(X) & \cdots & \nabla \nabla_{X(:,1)} g_N(X) \\ \nabla \nabla_{X(:,2)} g_1(X) & \nabla \nabla_{X(:,2)} g_2(X) & \cdots & \nabla \nabla_{X(:,2)} g_N(X) \\ \ddots & \ddots & \ddots & \ddots \\ \nabla \nabla_{X(:,L)} g_1(X) & \nabla \nabla_{X(:,L)} g_2(X) & \cdots & \nabla \nabla_{X(:,L)} g_N(X) \end{bmatrix} \quad (1944)$$

$$= [\nabla^2 g_1(X) \ \nabla^2 g_2(X) \ \cdots \ \nabla^2 g_N(X)] \in \mathbb{R}^{K \times N \times L \times K \times L}$$

The gradient of matrix-valued function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$  on matrix domain has a four-dimensional representation called *quartix* (*fourth-order tensor*)

$$\nabla g(X) \triangleq \begin{bmatrix} \nabla g_{11}(X) & \nabla g_{12}(X) & \cdots & \nabla g_{1N}(X) \\ \nabla g_{21}(X) & \nabla g_{22}(X) & \cdots & \nabla g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla g_{M1}(X) & \nabla g_{M2}(X) & \cdots & \nabla g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L} \quad (1945)$$

while the second-order gradient has a six-dimensional representation

$$\nabla^2 g(X) \triangleq \begin{bmatrix} \nabla^2 g_{11}(X) & \nabla^2 g_{12}(X) & \cdots & \nabla^2 g_{1N}(X) \\ \nabla^2 g_{21}(X) & \nabla^2 g_{22}(X) & \cdots & \nabla^2 g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla^2 g_{M1}(X) & \nabla^2 g_{M2}(X) & \cdots & \nabla^2 g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L \times K \times L} \quad (1946)$$

and so on.

### D.1.2 Product rules for matrix-functions

Given dimensionally compatible matrix-valued functions of matrix variable  $f(X)$  and  $g(X)$

$$\nabla_X(f(X)^T g(X)) = \nabla_X(f) g + \nabla_X(g) f \quad (1947)$$

while [56, §8.3] [352]

$$\nabla_X \text{tr}(f(X)^T g(X)) = \nabla_X \left( \text{tr}(f(X)^T g(Z)) + \text{tr}(g(X) f(Z)^T) \right) \Big|_{Z \leftarrow X} \quad (1948)$$

These expressions implicitly apply as well to scalar-, vector-, or matrix-valued functions of scalar, vector, or matrix arguments.

#### D.1.2.0.1 Example. Cubix.

Suppose  $f(X) : \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}^2 = X^T a$  and  $g(X) : \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}^2 = X b$ . We wish to find

$$\nabla_X(f(X)^T g(X)) = \nabla_X a^T X^2 b \quad (1949)$$

using the product rule. Formula (1947) calls for

$$\nabla_X a^T X^2 b = \nabla_X(X^T a) X b + \nabla_X(X b) X^T a \quad (1950)$$

Consider the first of the two terms:

$$\begin{aligned} \nabla_X(f) g &= \nabla_X(X^T a) X b \\ &= [\nabla(X^T a)_1 \quad \nabla(X^T a)_2] X b \end{aligned} \quad (1951)$$

The gradient of  $X^T a$  forms a cubix in  $\mathbb{R}^{2 \times 2 \times 2}$ ; a.k.a, *third-order tensor*.

$$\nabla_X(X^T a) X b = \left[ \begin{array}{ccc|cc} \frac{\partial(X^T a)_1}{\partial X_{11}} & \cdots & \cdots & \frac{\partial(X^T a)_2}{\partial X_{11}} & \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{\partial(X^T a)_1}{\partial X_{12}} & \cdots & \cdots & \frac{\partial(X^T a)_2}{\partial X_{12}} & \\ \hline \frac{\partial(X^T a)_1}{\partial X_{21}} & \cdots & \cdots & \frac{\partial(X^T a)_2}{\partial X_{21}} & \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{\partial(X^T a)_1}{\partial X_{22}} & \cdots & \cdots & \frac{\partial(X^T a)_2}{\partial X_{22}} & \end{array} \right] \begin{bmatrix} (X b)_1 \\ (X b)_2 \end{bmatrix} \in \mathbb{R}^{2 \times 1 \times 2} \quad (1952)$$

Because gradient of the product (1949) requires total change with respect to change in each entry of matrix  $X$ , the  $X b$  vector must make an inner product with each vector in that second dimension of the cubix indicated by dotted line segments;

$$\begin{aligned} \nabla_X(X^T a) X b &= \begin{bmatrix} a_1 & 0 & a_1 \\ a_2 & 0 & a_2 \\ 0 & a_1 & a_2 \end{bmatrix} \begin{bmatrix} b_1 X_{11} + b_2 X_{12} \\ b_1 X_{21} + b_2 X_{22} \end{bmatrix} \in \mathbb{R}^{2 \times 1 \times 2} \\ &= \begin{bmatrix} a_1(b_1 X_{11} + b_2 X_{12}) & a_1(b_1 X_{21} + b_2 X_{22}) \\ a_2(b_1 X_{11} + b_2 X_{12}) & a_2(b_1 X_{21} + b_2 X_{22}) \end{bmatrix} \in \mathbb{R}^{2 \times 2} \\ &= ab^T X^T \end{aligned} \quad (1953)$$

where the cubix appears as a complete  $2 \times 2 \times 2$  matrix. In like manner for the second term  $\nabla_X(g) f$

$$\begin{aligned}\nabla_X(Xb) X^T a &= \begin{bmatrix} b_1 & 0 \\ 0 & b_2 \\ 0 & 0 \\ 0 & b_1 \\ 0 & b_2 \end{bmatrix} \begin{bmatrix} X_{11}a_1 + X_{21}a_2 \\ X_{12}a_1 + X_{22}a_2 \end{bmatrix} \in \mathbb{R}^{2 \times 1 \times 2} \\ &= X^T a b^T \in \mathbb{R}^{2 \times 2}\end{aligned}\quad (1954)$$

The solution

$$\nabla_X a^T X^2 b = ab^T X^T + X^T a b^T \quad (1955)$$

can be found from Table D.2.1 or verified using (1948).  $\square$

### D.1.2.1 Kronecker product

A partial remedy for venturing into *hyperdimensional* matrix representations, such as the cubix or quartix, is to first vectorize matrices as in (37). This device gives rise to the Kronecker product of matrices  $\otimes$ ; a.k.a, *tensor product* (`kron()` in Matlab). Although its definition sees reversal in the literature, [363, §2.1] Kronecker product is not commutative ( $B \otimes A \neq A \otimes B$ ). We adopt the definition: for  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{p \times q}$

$$B \otimes A \triangleq \begin{bmatrix} B_{11}A & B_{12}A & \cdots & B_{1q}A \\ B_{21}A & B_{22}A & \cdots & B_{2q}A \\ \vdots & \vdots & & \vdots \\ B_{p1}A & B_{p2}A & \cdots & B_{pq}A \end{bmatrix} \in \mathbb{R}^{pm \times qn} \quad (1956)$$

for which  $A \otimes 1 = 1 \otimes A = A$  (real unity acts like Identity).

One advantage to vectorization is existence of the traditional two-dimensional matrix representation (*second-order tensor*) for the second-order gradient of a real function with respect to a vectorized matrix. From §A.1.1 no.36 (§D.2.1) for square  $A, B \in \mathbb{R}^{n \times n}$ , for example [190, §5.2] [13, §3]

$$\nabla_{\text{vec } X}^2 \text{tr}(AXBX^T) = \nabla_{\text{vec } X}^2 \text{vec}(X)^T (B^T \otimes A) \text{vec } X = B \otimes A^T + B^T \otimes A \in \mathbb{R}^{n^2 \times n^2} \quad (1957)$$

To disadvantage is a large new but known set of algebraic rules (§A.1.1) and the fact that its mere use does not generally guarantee two-dimensional matrix representation of gradients.

Another application of the Kronecker product is to reverse order of appearance in a matrix product: Suppose we wish to weight the columns of a matrix  $S \in \mathbb{R}^{M \times N}$ , for example, by respective entries  $w_i$  from the main diagonal in

$$W \triangleq \begin{bmatrix} w_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & w_N \end{bmatrix} \in \mathbb{S}^N \quad (1958)$$

A conventional means for accomplishing column weighting is to multiply  $S$  by diagonal matrix  $W$  on the right-hand side:

$$SW = S \begin{bmatrix} w_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & w_N \end{bmatrix} = [S(:, 1)w_1 \quad \cdots \quad S(:, N)w_N] \in \mathbb{R}^{M \times N} \quad (1959)$$

To reverse product order such that diagonal matrix  $W$  instead appears to the left of  $S$ : for  $I \in \mathbb{S}^M$  (Law)

$$SW = (\delta(W)^T \otimes I) \begin{bmatrix} S(:, 1) & 0 & & \mathbf{0} \\ 0 & S(:, 2) & \ddots & \\ & \ddots & \ddots & 0 \\ \mathbf{0} & & 0 & S(:, N) \end{bmatrix} \in \mathbb{R}^{M \times N} \quad (1960)$$

To instead weight the rows of  $S$  via diagonal matrix  $W \in \mathbb{S}^M$ , for  $I \in \mathbb{S}^N$

$$WS = \begin{bmatrix} S(1, :) & 0 & & \mathbf{0} \\ 0 & S(2, :) & \ddots & \\ & \ddots & \ddots & 0 \\ \mathbf{0} & & 0 & S(M, :) \end{bmatrix} (\delta(W) \otimes I) \in \mathbb{R}^{M \times N} \quad (1961)$$

### D.1.2.2 Hadamard product

For any matrices of like size,  $S, Y \in \mathbb{R}^{M \times N}$ , Hadamard's product  $\circ$  denotes simple multiplication of corresponding entries ( $.*$  in Matlab). It is possible to convert Hadamard product into a standard product of matrices:

$$S \circ Y = [\delta(Y(:, 1)) \cdots \delta(Y(:, N))] \begin{bmatrix} S(:, 1) & 0 & & \mathbf{0} \\ 0 & S(:, 2) & \ddots & \\ & \ddots & \ddots & 0 \\ \mathbf{0} & & 0 & S(:, N) \end{bmatrix} \in \mathbb{R}^{M \times N} \quad (1962)$$

In the special case that  $S = s$  and  $Y = y$  are vectors in  $\mathbb{R}^M$

$$s \circ y = \delta(s)y \quad (1963)$$

$$\begin{aligned} s^T \otimes y &= ys^T \\ s \otimes y^T &= sy^T \end{aligned} \quad (1964)$$

### D.1.3 Chain rules for composite matrix-functions

Given dimensionally compatible matrix-valued functions of matrix variable  $f(X)$  and  $g(X)$  [387, §15.7]

$$\nabla_X g(f(X)^T) = \nabla_X f^T \nabla_f g \quad (1965)$$

$$\nabla_X^2 g(f(X)^T) = \nabla_X (\nabla_X f^T \nabla_f g) = \nabla_X^2 f \nabla_f g + \nabla_X f^T \nabla_f^2 g \nabla_X f \quad (1966)$$

#### D.1.3.1 Two arguments

$$\nabla_X g(f(X)^T, h(X)^T) = \nabla_X f^T \nabla_f g + \nabla_X h^T \nabla_h g \quad (1967)$$

##### D.1.3.1.1 Example. Chain rule for two arguments.

[43, §1.1]

$$g(f(x)^T, h(x)^T) = (f(x) + h(x))^T A (f(x) + h(x)) \quad (1968)$$

$$f(x) = \begin{bmatrix} x_1 \\ \varepsilon x_2 \end{bmatrix}, \quad h(x) = \begin{bmatrix} \varepsilon x_1 \\ x_2 \end{bmatrix} \quad (1969)$$

$$\nabla_x g(f(x)^T, h(x)^T) = \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \end{bmatrix}(A + A^T)(f + h) + \begin{bmatrix} \varepsilon & 0 \\ 0 & 1 \end{bmatrix}(A + A^T)(f + h) \quad (1970)$$

$$\nabla_x g(f(x)^T, h(x)^T) = \begin{bmatrix} 1 + \varepsilon & 0 \\ 0 & 1 + \varepsilon \end{bmatrix}(A + A^T) \left( \begin{bmatrix} x_1 \\ \varepsilon x_2 \end{bmatrix} + \begin{bmatrix} \varepsilon x_1 \\ x_2 \end{bmatrix} \right) \quad (1971)$$

$$\lim_{\varepsilon \rightarrow 0} \nabla_x g(f(x)^T, h(x)^T) = (A + A^T)x \quad (1972)$$

from Table D.2.1.  $\square$

These foregoing formulae remain correct when gradient produces hyperdimensional representation:

#### D.1.4 First directional derivative

Assume that a differentiable function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$  has continuous first- and second-order gradients  $\nabla g$  and  $\nabla^2 g$  over  $\text{dom } g$  which is an open set. We seek simple expressions for the first and second directional derivatives in direction  $Y \in \mathbb{R}^{K \times L}$ : respectively,  $\overset{\rightarrow}{dg} \in \mathbb{R}^{M \times N}$  and  $\overset{\rightarrow}{dg^2} \in \mathbb{R}^{M \times N}$ .

Assuming that the limit exists, we may state the partial derivative of the  $mn^{\text{th}}$  entry of  $g$  with respect to  $kl^{\text{th}}$  entry of  $X$ :

$$\frac{\partial g_{mn}(X)}{\partial X_{kl}} = \lim_{\Delta t \rightarrow 0} \frac{g_{mn}(X + \Delta t e_k e_l^T) - g_{mn}(X)}{\Delta t} \in \mathbb{R} \quad (1973)$$

where  $e_k$  is the  $k^{\text{th}}$  standard basis vector in  $\mathbb{R}^K$  while  $e_l$  is the  $l^{\text{th}}$  standard basis vector in  $\mathbb{R}^L$ . Total number of partial derivatives equals  $KLMN$  while the gradient is defined in their terms;  $mn^{\text{th}}$  entry of the gradient is

$$\nabla g_{mn}(X) = \begin{bmatrix} \frac{\partial g_{mn}(X)}{\partial X_{11}} & \frac{\partial g_{mn}(X)}{\partial X_{12}} & \dots & \frac{\partial g_{mn}(X)}{\partial X_{1L}} \\ \frac{\partial g_{mn}(X)}{\partial X_{21}} & \frac{\partial g_{mn}(X)}{\partial X_{22}} & \dots & \frac{\partial g_{mn}(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial g_{mn}(X)}{\partial X_{K1}} & \frac{\partial g_{mn}(X)}{\partial X_{K2}} & \dots & \frac{\partial g_{mn}(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L} \quad (1974)$$

while the gradient is a quartix

$$\nabla g(X) = \begin{bmatrix} \nabla g_{11}(X) & \nabla g_{12}(X) & \dots & \nabla g_{1N}(X) \\ \nabla g_{21}(X) & \nabla g_{22}(X) & \dots & \nabla g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla g_{M1}(X) & \nabla g_{M2}(X) & \dots & \nabla g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L} \quad (1975)$$

By simply rotating our perspective of a four-dimensional representation of gradient matrix, we find one of three useful transpositions of this quartix (connote  $T_1$ ):

$$\nabla g(X)^{T_1} = \begin{bmatrix} \frac{\partial g(X)}{\partial X_{11}} & \frac{\partial g(X)}{\partial X_{12}} & \dots & \frac{\partial g(X)}{\partial X_{1L}} \\ \frac{\partial g(X)}{\partial X_{21}} & \frac{\partial g(X)}{\partial X_{22}} & \dots & \frac{\partial g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial g(X)}{\partial X_{K1}} & \frac{\partial g(X)}{\partial X_{K2}} & \dots & \frac{\partial g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times M \times N} \quad (1976)$$

When a limit for  $\Delta t \in \mathbb{R}$  exists, it is easy to show by substitution of variables in (1973)

$$\frac{\partial g_{mn}(X)}{\partial X_{kl}} Y_{kl} = \lim_{\Delta t \rightarrow 0} \frac{g_{mn}(X + \Delta t Y_{kl} e_k e_l^T) - g_{mn}(X)}{\Delta t} \in \mathbb{R} \quad (1977)$$

which may be interpreted as the change in  $g_{mn}$  at  $X$  when the change in  $X_{kl}$  is equal to  $Y_{kl}$  the  $kl^{\text{th}}$  entry of any  $Y \in \mathbb{R}^{K \times L}$ . Because the total change in  $g_{mn}(X)$  due to  $Y$  is the sum of change with respect to each and every  $X_{kl}$ , the  $mn^{\text{th}}$  entry of the directional derivative is the corresponding total differential [387, §15.8]

$$dg_{mn}(X)|_{dX \rightarrow Y} = \sum_{k,l} \frac{\partial g_{mn}(X)}{\partial X_{kl}} Y_{kl} = \text{tr}(\nabla g_{mn}(X)^T Y) \quad (1978)$$

$$= \sum_{k,l} \lim_{\Delta t \rightarrow 0} \frac{g_{mn}(X + \Delta t Y_{kl} e_k e_l^T) - g_{mn}(X)}{\Delta t} \quad (1979)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{g_{mn}(X + \Delta t Y) - g_{mn}(X)}{\Delta t} \quad (1980)$$

$$= \frac{d}{dt} \Big|_{t=0} g_{mn}(X + t Y) \quad (1981)$$

where  $t \in \mathbb{R}$ . Assuming finite  $Y$ , equation (1980) is called the *Gâteaux differential* [42, App.A.5] [225, §D.2.1] [399, §5.28] whose existence is implied by existence of the *Fréchet differential* (the sum in (1978)). [280, §7.2] Each may be understood as the change in  $g_{mn}$  at  $X$  when the change in  $X$  is equal in magnitude and direction to  $Y$ .<sup>D.2</sup> Hence the directional derivative,

$$\begin{aligned} \overset{\rightarrow}{dg}(X) &\triangleq \left[ \begin{array}{cccc} dg_{11}(X) & dg_{12}(X) & \cdots & dg_{1N}(X) \\ dg_{21}(X) & dg_{22}(X) & \cdots & dg_{2N}(X) \\ \vdots & \vdots & & \vdots \\ dg_{M1}(X) & dg_{M2}(X) & \cdots & dg_{MN}(X) \end{array} \right] \Big|_{dX \rightarrow Y} \in \mathbb{R}^{M \times N} \\ &= \left[ \begin{array}{cccc} \text{tr}(\nabla g_{11}(X)^T Y) & \text{tr}(\nabla g_{12}(X)^T Y) & \cdots & \text{tr}(\nabla g_{1N}(X)^T Y) \\ \text{tr}(\nabla g_{21}(X)^T Y) & \text{tr}(\nabla g_{22}(X)^T Y) & \cdots & \text{tr}(\nabla g_{2N}(X)^T Y) \\ \vdots & \vdots & & \vdots \\ \text{tr}(\nabla g_{M1}(X)^T Y) & \text{tr}(\nabla g_{M2}(X)^T Y) & \cdots & \text{tr}(\nabla g_{MN}(X)^T Y) \end{array} \right] \quad (1982) \\ &= \left[ \begin{array}{cccc} \sum_{k,l} \frac{\partial g_{11}(X)}{\partial X_{kl}} Y_{kl} & \sum_{k,l} \frac{\partial g_{12}(X)}{\partial X_{kl}} Y_{kl} & \cdots & \sum_{k,l} \frac{\partial g_{1N}(X)}{\partial X_{kl}} Y_{kl} \\ \sum_{k,l} \frac{\partial g_{21}(X)}{\partial X_{kl}} Y_{kl} & \sum_{k,l} \frac{\partial g_{22}(X)}{\partial X_{kl}} Y_{kl} & \cdots & \sum_{k,l} \frac{\partial g_{2N}(X)}{\partial X_{kl}} Y_{kl} \\ \vdots & \vdots & & \vdots \\ \sum_{k,l} \frac{\partial g_{M1}(X)}{\partial X_{kl}} Y_{kl} & \sum_{k,l} \frac{\partial g_{M2}(X)}{\partial X_{kl}} Y_{kl} & \cdots & \sum_{k,l} \frac{\partial g_{MN}(X)}{\partial X_{kl}} Y_{kl} \end{array} \right] \end{aligned}$$

from which it follows

$$\overset{\rightarrow}{dg}(X) = \sum_{k,l} \frac{\partial g(X)}{\partial X_{kl}} Y_{kl} \quad (1983)$$

---

<sup>D.2</sup>Although  $Y$  is a matrix, we may regard it as a vector in  $\mathbb{R}^{KL}$ .

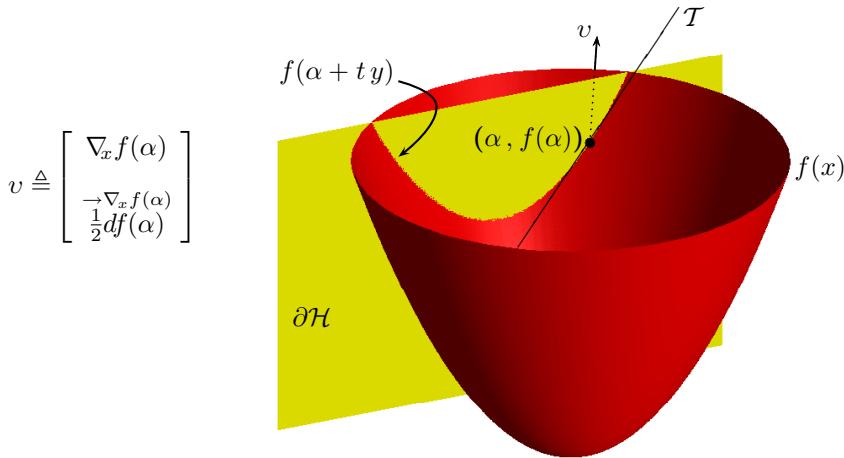


Figure 186: Strictly convex quadratic bowl in  $\mathbb{R}^2 \times \mathbb{R}$ ;  $f(x) = x^T x : \mathbb{R}^2 \rightarrow \mathbb{R}$  versus  $x$  on some open disc in  $\mathbb{R}^2$ . Plane slice  $\partial\mathcal{H}$  is perpendicular to function domain. Slice intersection with domain connotes bidirectional vector  $y$ . Slope of tangent line  $T$  at point  $(\alpha, f(\alpha))$  is value of directional derivative  $\nabla_x f(\alpha)^T y$  (2010) at  $\alpha$  in slice direction  $y$ . Negative gradient  $-\nabla_x f(x) \in \mathbb{R}^2$  is direction of *steepest descent*. [387, §15.6] [173] When vector  $v \in \mathbb{R}^3$  entry  $v_3$  is half directional derivative in gradient direction at  $\alpha$  and when  $\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \nabla_x f(\alpha)$ , then  $-v$  points directly toward bowl bottom.

Yet for all  $X \in \text{dom } g$ , any  $Y \in \mathbb{R}^{K \times L}$ , and some open interval of  $t \in \mathbb{R}$

$$g(X + tY) = g(X) + t \overset{\rightarrow}{dg}(X) + O(t^2) \quad (1984)$$

which is the first-order multidimensional Taylor series expansion about  $X$ . [387, §18.4] [173, §2.3.4] Differentiation with respect to  $t$  and subsequent  $t$ -zeroing isolates the second term of expansion. Thus differentiating and zeroing  $g(X + tY)$  in  $t$  is an operation equivalent to individually differentiating and zeroing every entry  $g_{mn}(X + tY)$  as in (1981). So the directional derivative of  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$  in any direction  $Y \in \mathbb{R}^{K \times L}$  evaluated at  $X \in \text{dom } g$  becomes

$$\overset{\rightarrow}{dg}(X) = \left. \frac{d}{dt} \right|_{t=0} g(X + tY) \in \mathbb{R}^{M \times N} \quad (1985)$$

[309, §2.1, §5.4.5] [35, §6.3.1] which is simplest. In case of a real function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$

$$\overset{\rightarrow}{dg}(X) = \text{tr}(\nabla g(X)^T Y) \quad (2007)$$

In case  $g(X) : \mathbb{R}^K \rightarrow \mathbb{R}$

$$\overset{\rightarrow}{dg}(X) = \nabla g(X)^T Y \quad (2010)$$

Unlike gradient, directional derivative does not expand dimension; directional derivative (1985) retains the dimensions of  $g$ . The derivative with respect to  $t$  makes the directional derivative resemble ordinary calculus (§D.2); e.g., when  $g(X)$  is linear,  $\overset{\rightarrow}{dg}(X) = g(Y)$ . [280, §7.2]

#### D.1.4.1 Interpretation of directional derivative

In the case of any differentiable real function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$ , the directional derivative of  $g(X)$  at  $X$  in any direction  $Y$  yields the slope of  $g$  along the line  $\{X + tY \mid t \in \mathbb{R}\}$  through its domain evaluated at  $t = 0$ . For higher-dimensional functions, by (1982), this slope interpretation can be applied to each entry of the directional derivative.

Figure 186, for example, shows a plane slice of a real convex bowl-shaped function  $f(x)$  along a line  $\{\alpha + ty \mid t \in \mathbb{R}\}$  through its domain. The slice reveals a one-dimensional real function of  $t$ ;  $f(\alpha + ty)$ . The directional derivative at  $x = \alpha$  in direction  $y$  is the slope of  $f(\alpha + ty)$  with respect to  $t$  at  $t = 0$ . In the case of a real function having vector argument  $h(X) : \mathbb{R}^K \rightarrow \mathbb{R}$ , its directional derivative in the normalized direction of its gradient is the gradient magnitude. (2010) For a real function of real variable, the directional derivative evaluated at any point in the function domain is just the slope of that function there scaled by the real direction. (*confer* §3.6)

Directional derivative generalizes our one-dimensional notion of derivative to a multidimensional domain. When direction  $Y$  coincides with a member of the standard Cartesian basis  $e_k e_l^T$  (61), then a single partial derivative  $\partial g(X)/\partial X_{kl}$  is obtained from directional derivative (1983); such is each entry of gradient  $\nabla g(X)$  in equalities (2007) and (2010), for example.

**D.1.4.1.1 Theorem.** *Directional derivative optimality condition.* [280, §7.4]  
Suppose  $f(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$  is minimized on convex set  $\mathcal{C} \subseteq \mathbb{R}^{K \times L}$  by  $X^*$ , and the directional derivative of  $f$  exists there. Then for all  $X \in \mathcal{C}$

$$\stackrel{\rightarrow}{df}(X) \geq 0 \quad (1986)$$

◊

#### D.1.4.1.2 Example. Simple bowl.

Bowl function (Figure 186)

$$f(x) : \mathbb{R}^K \rightarrow \mathbb{R} \triangleq (x - a)^T (x - a) - b \quad (1987)$$

has function offset  $-b \in \mathbb{R}$ , axis of revolution at  $x = a$ , and positive definite Hessian (1936) everywhere in its domain (an open *hyperdisc* in  $\mathbb{R}^K$ ); *id est*, strictly convex quadratic  $f(x)$  has unique global minimum equal to  $-b$  at  $x = a$ . A vector  $-v$  based anywhere in  $\text{dom } f \times \mathbb{R}$  pointing toward the unique bowl-bottom is specified:

$$v \propto \begin{bmatrix} x - a \\ f(x) + b \end{bmatrix} \in \mathbb{R}^K \times \mathbb{R} \quad (1988)$$

Such a vector is

$$v = \begin{bmatrix} \nabla_x f(x) \\ \stackrel{\rightarrow}{\nabla_x f(x)} \\ \frac{1}{2} df(x) \end{bmatrix} \quad (1989)$$

since the gradient is

$$\nabla_x f(x) = 2(x - a) \quad (1990)$$

and the directional derivative in direction of the gradient is (2010)

$$\stackrel{\rightarrow}{df}(x) = \nabla_x f(x)^T \nabla_x f(x) = 4(x - a)^T (x - a) = 4(f(x) + b) \quad (1991)$$

□

### D.1.5 Second directional derivative

By similar argument, it so happens: the second directional derivative is equally simple. Given  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$  on open domain,

$$\nabla \frac{\partial g_{mn}(X)}{\partial X_{kl}} = \frac{\partial \nabla g_{mn}(X)}{\partial X_{kl}} = \begin{bmatrix} \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{11}} & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{12}} & \dots & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{1L}} \\ \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{21}} & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{22}} & \dots & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{K1}} & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{K2}} & \dots & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L} \quad (1992)$$

$$\nabla^2 g_{mn}(X) = \begin{bmatrix} \nabla \frac{\partial g_{mn}(X)}{\partial X_{11}} & \nabla \frac{\partial g_{mn}(X)}{\partial X_{12}} & \dots & \nabla \frac{\partial g_{mn}(X)}{\partial X_{1L}} \\ \nabla \frac{\partial g_{mn}(X)}{\partial X_{21}} & \nabla \frac{\partial g_{mn}(X)}{\partial X_{22}} & \dots & \nabla \frac{\partial g_{mn}(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \nabla \frac{\partial g_{mn}(X)}{\partial X_{K1}} & \nabla \frac{\partial g_{mn}(X)}{\partial X_{K2}} & \dots & \nabla \frac{\partial g_{mn}(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times K \times L} \quad (1993)$$

$$= \begin{bmatrix} \frac{\partial \nabla g_{mn}(X)}{\partial X_{11}} & \frac{\partial \nabla g_{mn}(X)}{\partial X_{12}} & \dots & \frac{\partial \nabla g_{mn}(X)}{\partial X_{1L}} \\ \frac{\partial \nabla g_{mn}(X)}{\partial X_{21}} & \frac{\partial \nabla g_{mn}(X)}{\partial X_{22}} & \dots & \frac{\partial \nabla g_{mn}(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial \nabla g_{mn}(X)}{\partial X_{K1}} & \frac{\partial \nabla g_{mn}(X)}{\partial X_{K2}} & \dots & \frac{\partial \nabla g_{mn}(X)}{\partial X_{KL}} \end{bmatrix}$$

Rotating our perspective, we get several views of the second-order gradient:

$$\nabla^2 g(X) = \begin{bmatrix} \nabla^2 g_{11}(X) & \nabla^2 g_{12}(X) & \dots & \nabla^2 g_{1N}(X) \\ \nabla^2 g_{21}(X) & \nabla^2 g_{22}(X) & \dots & \nabla^2 g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla^2 g_{M1}(X) & \nabla^2 g_{M2}(X) & \dots & \nabla^2 g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L \times K \times L} \quad (1994)$$

$$\nabla^2 g(X)^T_1 = \begin{bmatrix} \nabla \frac{\partial g(X)}{\partial X_{11}} & \nabla \frac{\partial g(X)}{\partial X_{12}} & \dots & \nabla \frac{\partial g(X)}{\partial X_{1L}} \\ \nabla \frac{\partial g(X)}{\partial X_{21}} & \nabla \frac{\partial g(X)}{\partial X_{22}} & \dots & \nabla \frac{\partial g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \nabla \frac{\partial g(X)}{\partial X_{K1}} & \nabla \frac{\partial g(X)}{\partial X_{K2}} & \dots & \nabla \frac{\partial g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times M \times N \times K \times L} \quad (1995)$$

$$\nabla^2 g(X)^T_2 = \begin{bmatrix} \frac{\partial \nabla g(X)}{\partial X_{11}} & \frac{\partial \nabla g(X)}{\partial X_{12}} & \dots & \frac{\partial \nabla g(X)}{\partial X_{1L}} \\ \frac{\partial \nabla g(X)}{\partial X_{21}} & \frac{\partial \nabla g(X)}{\partial X_{22}} & \dots & \frac{\partial \nabla g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial \nabla g(X)}{\partial X_{K1}} & \frac{\partial \nabla g(X)}{\partial X_{K2}} & \dots & \frac{\partial \nabla g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times K \times L \times M \times N} \quad (1996)$$

Assuming the limits to exist, we may state the partial derivative of the  $mn^{\text{th}}$  entry of  $g$  with respect to  $kl^{\text{th}}$  and  $ij^{\text{th}}$  entries of  $X$ ;

$$\begin{aligned} \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{ij}} &= \frac{\partial}{\partial X_{ij}} \left( \frac{\partial g_{mn}(X)}{\partial X_{kl}} \right) = \lim_{\Delta t \rightarrow 0} \frac{\partial g_{mn}(X + \Delta t e_k e_l^T) - \partial g_{mn}(X)}{\partial X_{ij} \Delta t} \\ &= \lim_{\Delta \tau, \Delta t \rightarrow 0} \frac{(g_{mn}(X + \Delta t e_k e_l^T + \Delta \tau e_i e_j^T) - g_{mn}(X + \Delta t e_k e_l^T)) - (g_{mn}(X + \Delta \tau e_i e_j^T) - g_{mn}(X))}{\Delta \tau \Delta t} \end{aligned} \quad (1997)$$

Differentiating (1977) and then scaling by  $Y_{ij}$

$$\begin{aligned} \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} &= \lim_{\Delta t \rightarrow 0} \frac{\partial g_{mn}(X + \Delta t Y_{kl} e_k e_l^T) - \partial g_{mn}(X)}{\partial X_{ij} \Delta t} Y_{ij} \\ &= \lim_{\Delta \tau, \Delta t \rightarrow 0} \frac{(g_{mn}(X + \Delta t Y_{kl} e_k e_l^T + \Delta \tau Y_{ij} e_i e_j^T) - g_{mn}(X + \Delta t Y_{kl} e_k e_l^T)) - (g_{mn}(X + \Delta \tau Y_{ij} e_i e_j^T) - g_{mn}(X))}{\Delta \tau \Delta t} \end{aligned} \quad (1998)$$

which can be proved by substitution of variables in (1997). The  $mn^{\text{th}}$  second-order total differential due to any  $Y \in \mathbb{R}^{K \times L}$  is

$$d^2 g_{mn}(X)|_{dX \rightarrow Y} = \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{mn}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} = \text{tr}(\nabla_X \text{tr}(\nabla g_{mn}(X)^T Y)^T Y) \quad (1999)$$

$$= \sum_{i,j} \lim_{\Delta t \rightarrow 0} \frac{\partial g_{mn}(X + \Delta t Y) - \partial g_{mn}(X)}{\partial X_{ij} \Delta t} Y_{ij} \quad (2000)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{g_{mn}(X + 2\Delta t Y) - 2g_{mn}(X + \Delta t Y) + g_{mn}(X)}{\Delta t^2} \quad (2001)$$

$$= \left. \frac{d^2}{dt^2} \right|_{t=0} g_{mn}(X + t Y) \quad (2002)$$

Hence the second directional derivative,

$$\begin{aligned} \stackrel{\rightarrow Y}{dg^2}(X) &\triangleq \left[ \begin{array}{cccc} d^2 g_{11}(X) & d^2 g_{12}(X) & \cdots & d^2 g_{1N}(X) \\ d^2 g_{21}(X) & d^2 g_{22}(X) & \cdots & d^2 g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ d^2 g_{M1}(X) & d^2 g_{M2}(X) & \cdots & d^2 g_{MN}(X) \end{array} \right] \Big|_{dX \rightarrow Y} \in \mathbb{R}^{M \times N} \\ &= \left[ \begin{array}{cccc} \text{tr}(\nabla \text{tr}(\nabla g_{11}(X)^T Y)^T Y) & \text{tr}(\nabla \text{tr}(\nabla g_{12}(X)^T Y)^T Y) & \cdots & \text{tr}(\nabla \text{tr}(\nabla g_{1N}(X)^T Y)^T Y) \\ \text{tr}(\nabla \text{tr}(\nabla g_{21}(X)^T Y)^T Y) & \text{tr}(\nabla \text{tr}(\nabla g_{22}(X)^T Y)^T Y) & \cdots & \text{tr}(\nabla \text{tr}(\nabla g_{2N}(X)^T Y)^T Y) \\ \vdots & \vdots & & \vdots \\ \text{tr}(\nabla \text{tr}(\nabla g_{M1}(X)^T Y)^T Y) & \text{tr}(\nabla \text{tr}(\nabla g_{M2}(X)^T Y)^T Y) & \cdots & \text{tr}(\nabla \text{tr}(\nabla g_{MN}(X)^T Y)^T Y) \end{array} \right] \\ &= \left[ \begin{array}{cccc} \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{11}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} & \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{12}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} & \cdots & \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{1N}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} \\ \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{21}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} & \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{22}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} & \cdots & \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{2N}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} \\ \vdots & \vdots & & \vdots \\ \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{M1}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} & \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{M2}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} & \cdots & \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{MN}(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} \end{array} \right] \end{aligned} \quad (2003)$$

from which it follows

$$\stackrel{\rightarrow Y}{dg^2}(X) = \sum_{i,j} \sum_{k,l} \frac{\partial^2 g(X)}{\partial X_{kl} \partial X_{ij}} Y_{kl} Y_{ij} = \sum_{i,j} \frac{\partial}{\partial X_{ij}} \stackrel{\rightarrow Y}{dg}(X) Y_{ij} \quad (2004)$$

Yet for all  $X \in \text{dom } g$ , any  $Y \in \mathbb{R}^{K \times L}$ , and some open interval of  $t \in \mathbb{R}$

$$g(X + t Y) = g(X) + t \stackrel{\rightarrow Y}{dg}(X) + \frac{1}{2!} t^2 \stackrel{\rightarrow Y}{dg^2}(X) + O(t^3) \quad (2005)$$

which is the second-order multidimensional Taylor series expansion about  $X$ . [387, §18.4] [173, §2.3.4] Differentiating twice with respect to  $t$  and subsequent  $t$ -zeroing isolates the third term of the expansion. Thus differentiating and zeroing  $g(X+tY)$  in  $t$  is an operation equivalent to individually differentiating and zeroing every entry  $g_{mn}(X+tY)$  as in (2002). So the second directional derivative of  $g(X):\mathbb{R}^{K\times L}\rightarrow\mathbb{R}^{M\times N}$  becomes [309, §2.1, §5.4.5] [35, §6.3.1]

$$\stackrel{\rightarrow Y}{dg^2}(X) = \frac{d^2}{dt^2} \Big|_{t=0} g(X+tY) \in \mathbb{R}^{M\times N} \quad (2006)$$

which is again simplest. (*confer*(1985)) Directional derivative retains the dimensions of  $g$ .

### D.1.6 directional derivative expressions

In the case of a real function  $g(X):\mathbb{R}^{K\times L}\rightarrow\mathbb{R}$ , all its directional derivatives are in  $\mathbb{R}$ :

$$\stackrel{\rightarrow Y}{dg}(X) = \text{tr}(\nabla g(X)^T Y) \quad (2007)$$

$$\stackrel{\rightarrow Y}{dg^2}(X) = \text{tr}\left(\nabla_X \text{tr}(\nabla g(X)^T Y)^T Y\right) = \text{tr}\left(\nabla_X \stackrel{\rightarrow Y}{dg}(X)^T Y\right) \quad (2008)$$

$$\stackrel{\rightarrow Y}{dg^3}(X) = \text{tr}\left(\nabla_X \text{tr}\left(\nabla_X \text{tr}(\nabla g(X)^T Y)^T Y\right)^T Y\right) = \text{tr}\left(\nabla_X \stackrel{\rightarrow Y}{dg^2}(X)^T Y\right) \quad (2009)$$

In the case  $g(X):\mathbb{R}^K\rightarrow\mathbb{R}$  has vector argument, they further simplify:

$$\stackrel{\rightarrow Y}{dg}(X) = \nabla g(X)^T Y \quad (2010)$$

$$\stackrel{\rightarrow Y}{dg^2}(X) = Y^T \nabla^2 g(X) Y \quad (2011)$$

$$\stackrel{\rightarrow Y}{dg^3}(X) = \nabla_X (Y^T \nabla^2 g(X) Y)^T Y \quad (2012)$$

and so on.

### D.1.7 higher-order multidimensional Taylor series

Series expansions of the differentiable matrix-valued function  $g(X)$ , of matrix argument, were given earlier in (1984) and (2005). Assume that  $g(X)$  has continuous first-, second-, and third-order gradients over open set  $\text{dom } g$ . Then, for  $X \in \text{dom } g$  and any  $Y \in \mathbb{R}^{K\times L}$ , the Taylor series is expressed on some open interval of  $\mu \in \mathbb{R}$

$$g(X + \mu Y) = g(X) + \mu \stackrel{\rightarrow Y}{dg}(X) + \frac{1}{2!} \mu^2 \stackrel{\rightarrow Y}{dg^2}(X) + \frac{1}{3!} \mu^3 \stackrel{\rightarrow Y}{dg^3}(X) + O(\mu^4) \quad (2013)$$

or on some open interval of  $\|Y\|_2$

$$g(Y) = g(X) + \stackrel{\rightarrow Y-X}{dg}(X) + \frac{1}{2!} \stackrel{\rightarrow Y-X}{dg^2}(X) + \frac{1}{3!} \stackrel{\rightarrow Y-X}{dg^3}(X) + O(\|Y\|^4) \quad (2014)$$

which are third-order expansions about  $X$ . The *mean value theorem* from calculus is what insures finite order of the series. [387] [43, §1.1] [42, App.A.5] [225, §0.4] These somewhat unbelievable formulae<sup>D.3</sup> imply that a function can be determined over the whole of its domain by knowing its value and all its directional derivatives at a single point  $X$ .

---

<sup>D.3</sup> e.g., real continuous and differentiable function of real variable  $f(x)=e^{-1/x^2}$  has no Taylor series expansion about  $x=0$ , of any practical use, because each derivative equals 0 there.

**D.1.7.0.1 Example.** *Inverse-matrix function.*

Say  $g(Y) = Y^{-1}$ . From the table on page 560,

$$\stackrel{\rightarrow}{dg}(X) = \frac{d}{dt} \Big|_{t=0} g(X + tY) = -X^{-1}YX^{-1} \quad (2015)$$

$$\stackrel{\rightarrow}{dg^2}(X) = \frac{d^2}{dt^2} \Big|_{t=0} g(X + tY) = 2X^{-1}YX^{-1}YX^{-1} \quad (2016)$$

$$\stackrel{\rightarrow}{dg^3}(X) = \frac{d^3}{dt^3} \Big|_{t=0} g(X + tY) = -6X^{-1}YX^{-1}YX^{-1}YX^{-1} \quad (2017)$$

Let's find the Taylor series expansion of  $g$  about  $X = I$ : Since  $g(I) = I$ , for  $\|Y\|_2 < 1$  ( $\mu = 1$  in (2013))

$$g(I + Y) = (I + Y)^{-1} = I - Y + Y^2 - Y^3 + \dots \quad (2018)$$

If  $Y$  is small,  $(I + Y)^{-1} \approx I - Y$ .<sup>D.4</sup> Now we find Taylor series expansion about  $X$ :

$$g(X + Y) = (X + Y)^{-1} = X^{-1} - X^{-1}YX^{-1} + 2X^{-1}YX^{-1}YX^{-1} - \dots \quad (2019)$$

If  $Y$  is small,  $(X + Y)^{-1} \approx X^{-1} - X^{-1}YX^{-1}$ . □

**D.1.7.0.2 Exercise.** *log det.*

(confer [65, p.644])

Find the first three terms of a Taylor series expansion for  $\log \det Y$ . Specify an open interval over which the expansion holds in vicinity of  $X$ . ▼

## D.1.8 Correspondence of gradient to derivative

From the foregoing expressions for directional derivative, we derive a relationship between gradient with respect to matrix  $X$  and derivative with respect to real variable  $t$ :

### D.1.8.1 first-order

Removing evaluation at  $t = 0$  from (1985),<sup>D.5</sup> we find an expression for the directional derivative of  $g(X)$  in direction  $Y$  evaluated anywhere along a line  $\{X + tY \mid t \in \mathbb{R}\}$  intersecting  $\text{dom } g$

$$\stackrel{\rightarrow}{dg}(X + tY) = \frac{d}{dt}g(X + tY) \quad (2020)$$

In the general case  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$ , from (1978) and (1981) we find

$$\text{tr}(\nabla_X g_{mn}(X + tY)^T Y) = \frac{d}{dt}g_{mn}(X + tY) \quad (2021)$$

which is valid at  $t = 0$ , of course, when  $X \in \text{dom } g$ . In the important case of a real function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$ , from (2007) we have simply

$$\text{tr}(\nabla_X g(X + tY)^T Y) = \frac{d}{dt}g(X + tY) \quad (2022)$$

When  $g(X) : \mathbb{R}^K \rightarrow \mathbb{R}$  has vector argument,

$$\nabla_X g(X + tY)^T Y = \frac{d}{dt}g(X + tY) \quad (2023)$$

<sup>D.4</sup>Had we instead set  $g(Y) = (I + Y)^{-1}$ , then the equivalent expansion would have been about  $X = 0$ .

<sup>D.5</sup>Justified by replacing  $X$  with  $X + tY$  in (1978)-(1980); beginning,

$$dg_{mn}(X + tY)|_{dX=Y} = \sum_{k,l} \frac{\partial g_{mn}(X + tY)}{\partial X_{kl}} Y_{kl}$$

**D.1.8.1.1 Example.** *Gradient.*

$g(X) = w^T X^T X w$ ,  $X \in \mathbb{R}^{K \times L}$ ,  $w \in \mathbb{R}^L$ . Using the tables in §D.2,

$$\text{tr}(\nabla_X g(X + tY)^T Y) = \text{tr}(2ww^T(X^T + tY^T)Y) \quad (2024)$$

$$= 2w^T(X^T Y + tY^T Y)w \quad (2025)$$

Applying equivalence (2022),

$$\frac{d}{dt}g(X + tY) = \frac{d}{dt}w^T(X + tY)^T(X + tY)w \quad (2026)$$

$$= w^T(X^T Y + Y^T X + 2tY^T Y)w \quad (2027)$$

$$= 2w^T(X^T Y + tY^T Y)w \quad (2028)$$

which is the same as (2025). Hence, the equivalence is demonstrated.

It is easy to extract  $\nabla g(X)$  from (2028) knowing only (2022):

$$\begin{aligned} \text{tr}(\nabla_X g(X + tY)^T Y) &= 2w^T(X^T Y + tY^T Y)w \\ &= 2\text{tr}(ww^T(X^T + tY^T)Y) \\ \text{tr}(\nabla_X g(X)^T Y) &= 2\text{tr}(ww^T X^T Y) \\ &\Leftrightarrow \\ \nabla_X g(X) &= 2Xww^T \end{aligned} \quad (2029)$$

□

**D.1.8.2 second-order**

Likewise removing the evaluation at  $t = 0$  from (2006),

$$\overset{\rightarrow}{dg^2}(X + tY) = \frac{d^2}{dt^2}g(X + tY) \quad (2030)$$

we can find a similar relationship between second-order gradient and second derivative: In the general case  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$  from (1999) and (2002),

$$\text{tr}\left(\nabla_X \text{tr}(\nabla_X g_{mn}(X + tY)^T Y)^T Y\right) = \frac{d^2}{dt^2}g_{mn}(X + tY) \quad (2031)$$

In the case of a real function  $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$  we have, of course,

$$\text{tr}\left(\nabla_X \text{tr}(\nabla_X g(X + tY)^T Y)^T Y\right) = \frac{d^2}{dt^2}g(X + tY) \quad (2032)$$

From (2011), the simpler case, where real function  $g(X) : \mathbb{R}^K \rightarrow \mathbb{R}$  has vector argument,

$$Y^T \nabla_X^2 g(X + tY) Y = \frac{d^2}{dt^2}g(X + tY) \quad (2033)$$

**D.1.8.2.1 Example.** *Second-order gradient.*

We want to find  $\nabla^2 g(X) \in \mathbb{R}^{K \times K \times K \times K}$  given real function  $g(X) = \log \det X$  having domain  $\text{intr } \mathbb{S}_+^K$ . From the tables in §D.2,

$$h(X) \triangleq \nabla g(X) = X^{-1} \in \text{intr } \mathbb{S}_+^K \quad (2034)$$

so  $\nabla^2 g(X) = \nabla h(X)$ . By (2021) and (1984), for  $Y \in \mathbb{S}^K$

$$\text{tr}(\nabla h_{mn}(X)^T Y) = \frac{d}{dt} \Big|_{t=0} h_{mn}(X + tY) \quad (2035)$$

$$= \left( \frac{d}{dt} \Big|_{t=0} h(X + tY) \right)_{mn} \quad (2036)$$

$$= \left( \frac{d}{dt} \Big|_{t=0} (X + tY)^{-1} \right)_{mn} \quad (2037)$$

$$= -(X^{-1} Y X^{-1})_{mn} \quad (2038)$$

Setting  $Y$  to a member of  $\{e_k e_l^T \in \mathbb{R}^{K \times K} \mid k, l = 1 \dots K\}$ , and employing a property (39) of the trace function we find

$$\nabla^2 g(X)_{mnkl} = \text{tr}(\nabla h_{mn}(X)^T e_k e_l^T) = \nabla h_{mn}(X)_{kl} = -(X^{-1} e_k e_l^T X^{-1})_{mn} \quad (2039)$$

$$\nabla^2 g(X)_{kl} = \nabla h(X)_{kl} = -(X^{-1} e_k e_l^T X^{-1}) \in \mathbb{R}^{K \times K} \quad (2040)$$

□

From all these first- and second-order expressions, we may generate new ones by evaluating both sides at arbitrary  $t$  (in some open interval) but only after differentiation.

## D.2 Tables of gradients and derivatives

- Results may be validated numerically via *extrapolation*. [275, §5.4] [121] When algebraically proving results for symmetric matrices, it is critical to take gradients ignoring symmetry and to then substitute symmetric entries afterward. [190] [69]
- $i, j, k, \ell, K, L, m, n, M, N$  are integers, unless otherwise noted,  $a, b \in \mathbb{R}^n$ ,  $x, y \in \mathbb{R}^k$ ,  $A, B \in \mathbb{R}^{m \times n}$ ,  $X, Y \in \mathbb{R}^{K \times L}$ ,  $t, \mu \in \mathbb{R}$ .
- $x^\mu$  means  $\delta(\delta(x)^\mu)$  for  $\mu \in \mathbb{R}$ ; *id est*, entrywise vector exponentiation.  $\delta$  is the main-diagonal linear operator (1569).  $x^0 \triangleq \mathbf{1}$ ,  $X^0 \triangleq I$  if square.

- $\frac{d}{dx} \triangleq \begin{bmatrix} \frac{d}{dx_1} \\ \vdots \\ \frac{d}{dx_k} \end{bmatrix}$ ,  $\overset{\rightarrow}{y}(x)$ ,  $\overset{\rightarrow}{y^2}(x)$  (directional derivatives §D.1),  $\log x$ ,  $e^x$ ,  $|x|$ ,  $x/y$

(Hadamard quotient),  $\text{sgn } x$ ,  $\sqrt[k]{x}$  (entrywise square root), *etcetera*, are maps  $f : \mathbb{R}^k \rightarrow \mathbb{R}^k$  that maintain dimension; *e.g.* (§A.1.1)

$$\frac{d}{dx} x^{-1} \triangleq \nabla_x \mathbf{1}^T \delta(x)^{-1} \mathbf{1} \quad (2041)$$

- For  $A$  a scalar or square matrix, we have the Taylor series [84, §3.6]

$$e^A \triangleq \sum_{k=0}^{\infty} \frac{1}{k!} A^k \quad (2042)$$

Further, [368, §5.4]

$$e^A \succ 0 \quad \forall A \in \mathbb{S}^m \quad (2043)$$

- For all square  $A$  and integer  $k$

$$\det^k A = \det A^k \quad (2044)$$

### D.2.1 algebraic

$\nabla_x x = \nabla_x x^T = I \in \mathbb{R}^{k \times k}$	$\nabla_X X = \nabla_X X^T \triangleq I \in \mathbb{R}^{K \times L \times K \times L}$ (Identity)
$\nabla_x \mathbf{1}^T x = \nabla_x x^T \mathbf{1} = \mathbf{1} \in \mathbb{R}^k$	$\nabla_X \mathbf{1}^T X \mathbf{1} = \nabla_X \mathbf{1}^T X^T \mathbf{1} = \mathbf{1} \mathbf{1}^T \in \mathbb{R}^{K \times L}$
$\nabla_x (Ax - b) = A^T$	
$\nabla_x (x^T A - b^T) = A$	
$\nabla_x (Ax - b)^T (Ax - b) = 2A^T(Ax - b)$	
$\nabla_x^2 (Ax - b)^T (Ax - b) = 2A^T A$	
$\nabla_x \sqrt{(Ax - b)^T (Ax - b)} = A^T (Ax - b) / \ Ax - b\ _2 = \nabla_x \ Ax - b\ _2$	
$\nabla_x z^T  Ax - b  = A^T \delta(z) \operatorname{sgn}(Ax - b), z_i \neq 0 \Rightarrow (Ax - b)_i \neq 0$	
$\nabla_x \mathbf{1}^T  Ax - b  = A^T \operatorname{sgn}(Ax - b) = \nabla_x \ Ax - b\ _1$	
$\nabla_x \mathbf{1}^T f( Ax - b ) = A^T \delta\left(\frac{df(y)}{dy}\Big _{y= Ax-b }\right) \operatorname{sgn}(Ax - b)$	
$\nabla_x (x^T Ax + 2x^T By + y^T Cy) = (A + A^T)x + 2By$	$\nabla_X a^T X b = \nabla_X b^T X^T a = ab^T$
$\nabla_x (x + y)^T A (x + y) = (A + A^T)(x + y)$	$\nabla_X a^T X^2 b = X^T ab^T + ab^T X^T$
$\nabla_x^2 (x^T Ax + 2x^T By + y^T Cy) = A + A^T$	$\nabla_X a^T X^{-1} b = -X^{-T} ab^T X^{-T}$
$\nabla_x a^T x^T x b = 2xa^T b$	$\nabla_X (X^{-1})_{kl} = \frac{\partial X^{-1}}{\partial X_{kl}} = -X^{-1} e_k e_l^T X^{-1}, \quad \begin{matrix} \text{confer} \\ (1976) \\ (2040) \end{matrix}$
$\nabla_x a^T x x^T b = (ab^T + ba^T)x$	$\nabla_X a^T X^T X b = X(ab^T + ba^T)$
$\nabla_x a^T x^T x a = 2xa^T a$	$\nabla_X a^T X X^T b = (ab^T + ba^T)X$
$\nabla_x a^T x x^T a = 2aa^T x$	$\nabla_X a^T X^T X a = 2Xaa^T$
$\nabla_x a^T y x^T b = ba^T y$	$\nabla_X a^T X X^T a = 2aa^T X$
$\nabla_x a^T y^T x b = yb^T a$	$\nabla_X a^T Y X^T b = ba^T Y$
$\nabla_x a^T x y^T b = ab^T y$	$\nabla_X a^T Y^T X b = Yab^T$
$\nabla_x a^T x^T y b = ya^T b$	$\nabla_X a^T X Y^T b = ab^T Y$

**algebraic** continued

$$\frac{d}{dt}(X+tY) = Y$$

$$\begin{aligned}\frac{d}{dt}B^T(X+tY)^{-1}A &= -B^T(X+tY)^{-1}Y(X+tY)^{-1}A \\ \frac{d}{dt}B^T(X+tY)^{-T}A &= -B^T(X+tY)^{-T}Y^T(X+tY)^{-T}A \\ \frac{d}{dt}B^T(X+tY)^\mu A &= \dots, \quad -1 \leq \mu \leq 1, \quad X, Y \in \mathbb{S}_+^M\end{aligned}$$

$$\begin{aligned}\frac{d^2}{dt^2}B^T(X+tY)^{-1}A &= 2B^T(X+tY)^{-1}Y(X+tY)^{-1}Y(X+tY)^{-1}A \\ \frac{d^3}{dt^3}B^T(X+tY)^{-1}A &= -6B^T(X+tY)^{-1}Y(X+tY)^{-1}Y(X+tY)^{-1}Y(X+tY)^{-1}A \\ \frac{d}{dt}((X+tY)^TA(X+tY)) &= Y^TAX + X^TAY + 2tY^TAY \\ \frac{d^2}{dt^2}((X+tY)^TA(X+tY)) &= 2Y^TAY \\ \frac{d}{dt}((X+tY)^TA(X+tY))^{-1} &= -((X+tY)^TA(X+tY))^{-1}(Y^TAX + X^TAY + 2tY^TAY)((X+tY)^TA(X+tY))^{-1} \\ \frac{d}{dt}((X+tY)A(X+tY)) &= YAX + XAY + 2tYAY \\ \frac{d^2}{dt^2}((X+tY)A(X+tY)) &= 2YAY\end{aligned}$$

### D.2.2 trace Kronecker

$$\nabla_{\text{vec } X} \text{tr}(AXBX^T) = \nabla_{\text{vec } X} \text{vec}(X)^T(B^T \otimes A) \text{vec } X = (B \otimes A^T + B^T \otimes A) \text{vec } X$$

$$\nabla_{\text{vec } X}^2 \text{tr}(AXBX^T) = \nabla_{\text{vec } X}^2 \text{vec}(X)^T(B^T \otimes A) \text{vec } X = B \otimes A^T + B^T \otimes A \quad (1957)$$

### D.2.3 trace

$\nabla_x \mu x = \mu I$	$\nabla_X \operatorname{tr} \mu X = \nabla_X \mu \operatorname{tr} X = \mu I$
$\nabla_x \mathbf{1}^T \delta(x)^{-1} \mathbf{1} = \frac{d}{dx} x^{-1} = -x^{-2}$	$\nabla_X \operatorname{tr} X^{-1} = -X^{-2T}$
$\nabla_x \mathbf{1}^T \delta(x)^{-1} y = -\delta(x)^{-2} y$	$\nabla_X \operatorname{tr}(X^{-1} Y) = \nabla_X \operatorname{tr}(Y X^{-1}) = -X^{-T} Y^T X^{-T}$
$\frac{d}{dx} x^\mu = \mu x^{\mu-1}$	$\nabla_X \operatorname{tr} X^\mu = \mu X^{\mu-1}, \quad X \in \mathbb{S}^M$
	$\nabla_X \operatorname{tr} X^j = j X^{(j-1)T}$
$\nabla_x (b - a^T x)^{-1} = (b - a^T x)^{-2} a$	$\nabla_X \operatorname{tr}((B - AX)^{-1}) = ((B - AX)^{-2} A)^T$
$\nabla_x (b - a^T x)^\mu = -\mu (b - a^T x)^{\mu-1} a$	
$\nabla_x x^T y = \nabla_x y^T x = y$	$\nabla_X \operatorname{tr}(X^T Y) = \nabla_X \operatorname{tr}(Y X^T) = \nabla_X \operatorname{tr}(Y^T X) = \nabla_X \operatorname{tr}(XY^T) = Y$
$\nabla_x x^T x = 2x$	$\nabla_X \operatorname{tr}(X^T X) = \nabla_X \operatorname{tr}(XX^T) = 2X$
	$\nabla_X \operatorname{tr}(AXBX^T) = \nabla_X \operatorname{tr}(BX^T A) = A^T X B^T + A X B$
	$\nabla_X \operatorname{tr}(AXBX) = \nabla_X \operatorname{tr}(BX^T A) = A^T X^T B^T + B^T X^T A^T$
	$\nabla_X \operatorname{tr}(AXAXAXAX) = \nabla_X \operatorname{tr}(XAXAXAXA) = 4(AXAXAXA)^T$
	$\nabla_X \operatorname{tr}(AXAXAX) = \nabla_X \operatorname{tr}(XAXAXA) = 3(AXAXA)^T$
	$\nabla_X \operatorname{tr}(AXAX) = \nabla_X \operatorname{tr}(XAXA) = 2(AXA)^T$
	$\nabla_X \operatorname{tr}(AX) = \nabla_X \operatorname{tr}(XA) = A^T$
	$\nabla_X \operatorname{tr}(YX^k) = \nabla_X \operatorname{tr}(X^k Y) = \sum_{i=0}^{k-1} (X^i Y X^{k-1-i})^T$
	$\nabla_X \operatorname{tr}(X^T YY^T XX^T YY^T X) = 4YY^T XX^T YY^T X$
	$\nabla_X \operatorname{tr}(XY Y^T X^T XY Y^T X^T) = 4XY Y^T X^T XY Y^T$
	$\nabla_X \operatorname{tr}(Y^T XX^T Y) = \nabla_X \operatorname{tr}(X^T YY^T X) = 2YY^T X$
	$\nabla_X \operatorname{tr}(Y^T X^T XY) = \nabla_X \operatorname{tr}(XY Y^T X^T) = 2XY Y^T$
	$\nabla_X \operatorname{tr}((X + Y)^T (X + Y)) = 2(X + Y) = \nabla_X \ X + Y\ _F^2$
	$\nabla_X \operatorname{tr}((X + Y)(X + Y)) = 2(X + Y)^T$
	$\nabla_X \operatorname{tr}(A^T X B) = \nabla_X \operatorname{tr}(X^T A B^T) = AB^T$
	$\nabla_X \operatorname{tr}(A^T X^{-1} B) = \nabla_X \operatorname{tr}(X^{-T} A B^T) = -X^{-T} A B^T X^{-T}$
	$\nabla_X a^T X b = \nabla_X \operatorname{tr}(ba^T X) = \nabla_X \operatorname{tr}(Xba^T) = ab^T$
	$\nabla_X b^T X^T a = \nabla_X \operatorname{tr}(X^T ab^T) = \nabla_X \operatorname{tr}(ab^T X^T) = ab^T$
	$\nabla_X a^T X^{-1} b = \nabla_X \operatorname{tr}(X^{-T} ab^T) = -X^{-T} ab^T X^{-T}$
	$\nabla_X a^T X^\mu b = \dots$

**trace** continued

$\frac{d}{dt} \operatorname{tr} g(X + tY) = \operatorname{tr} \frac{d}{dt} g(X + tY)$ $\frac{d}{dt} \operatorname{tr}(X + tY) = \operatorname{tr} Y$ $\frac{d}{dt} \operatorname{tr}^j(X + tY) = j \operatorname{tr}^{j-1}(X + tY) \operatorname{tr} Y$ $\frac{d}{dt} \operatorname{tr}(X + tY)^j = j \operatorname{tr}((X + tY)^{j-1} Y) \quad (\forall j)$ $\frac{d}{dt} \operatorname{tr}((X + tY)Y) = \operatorname{tr} Y^2$ $\frac{d}{dt} \operatorname{tr}((X + tY)^k Y) = \frac{d}{dt} \operatorname{tr}(Y(X + tY)^k) = k \operatorname{tr}((X + tY)^{k-1} Y^2), \quad k \in \{0, 1, 2\}$ $\frac{d}{dt} \operatorname{tr}((X + tY)^k Y) = \frac{d}{dt} \operatorname{tr}(Y(X + tY)^k) = \operatorname{tr} \sum_{i=0}^{k-1} (X + tY)^i Y (X + tY)^{k-1-i} Y$ $\begin{aligned} \frac{d}{dt} \operatorname{tr}((X + tY)^{-1} Y) &= -\operatorname{tr}((X + tY)^{-1} Y (X + tY)^{-1} Y) \\ \frac{d}{dt} \operatorname{tr}(B^T (X + tY)^{-1} A) &= -\operatorname{tr}(B^T (X + tY)^{-1} Y (X + tY)^{-1} A) \\ \frac{d}{dt} \operatorname{tr}(B^T (X + tY)^{-T} A) &= -\operatorname{tr}(B^T (X + tY)^{-T} Y^T (X + tY)^{-T} A) \\ \frac{d}{dt} \operatorname{tr}(B^T (X + tY)^{-k} A) &= \dots, \quad k > 0 \\ \frac{d}{dt} \operatorname{tr}(B^T (X + tY)^{\mu} A) &= \dots, \quad -1 \leq \mu \leq 1, \quad X, Y \in \mathbb{S}_+^M \end{aligned}$ $\frac{d^2}{dt^2} \operatorname{tr}(B^T (X + tY)^{-1} A) = 2 \operatorname{tr}(B^T (X + tY)^{-1} Y (X + tY)^{-1} Y (X + tY)^{-1} A)$ $\begin{aligned} \frac{d}{dt} \operatorname{tr}((X + tY)^T A (X + tY)) &= \operatorname{tr}(Y^T A X + X^T A Y + 2t Y^T A Y) \\ \frac{d^2}{dt^2} \operatorname{tr}((X + tY)^T A (X + tY)) &= 2 \operatorname{tr}(Y^T A Y) \\ \frac{d}{dt} \operatorname{tr}\left(((X + tY)^T A (X + tY))^{-1}\right) &= -\operatorname{tr}\left(((X + tY)^T A (X + tY))^{-1} (Y^T A X + X^T A Y + 2t Y^T A Y) ((X + tY)^T A (X + tY))^{-1}\right) \\ \frac{d}{dt} \operatorname{tr}((X + tY) A (X + tY)) &= \operatorname{tr}(Y A X + X A Y + 2t Y A Y) \\ \frac{d^2}{dt^2} \operatorname{tr}((X + tY) A (X + tY)) &= 2 \operatorname{tr}(Y A Y) \end{aligned}$	[229, p.491]
--	--------------

#### D.2.4 logarithmic determinant

$x \succ 0$ ,  $\det X > 0$  on some neighborhood of  $X$ , and  $\det(X + tY) > 0$  on some open interval of  $t$ ; otherwise,  $\log(\cdot)$  would be discontinuous. [90, p.75]

$\frac{d}{dx} \log x = x^{-1}$	$\nabla_X \log \det X = X^{-T}$
	$\nabla_X^2 \log \det(X)_{kl} = \frac{\partial X^{-T}}{\partial X_{kl}} = -(X^{-1} e_k e_l^T X^{-1})^T$ , confer (1993)(2040)
$\frac{d}{dx} \log x^{-1} = -x^{-2}$	$\nabla_X \log \det X^{-1} = -X^{-T}$
$\frac{d}{dx} \log x^\mu = \mu x^{-1}$	$\nabla_X \log \det^\mu X = \mu X^{-T}$
	$\nabla_X \log \det X^\mu = \mu X^{-T}$
	$\nabla_X \log \det X^k = \nabla_X \log \det^k X = k X^{-T}$
	$\nabla_X \log \det^\mu(X + tY) = \mu(X + tY)^{-T}$
$\nabla_x \log(a^T x + b) = a \frac{1}{a^T x + b}$	$\nabla_X \log \det(AX + B) = A^T(AX + B)^{-T}$
	$\nabla_X \log \det(I \pm A^T X A) = \pm A(I \pm A^T X A)^{-T} A^T$
	$\nabla_X \log \det(X + tY)^k = \nabla_X \log \det^k(X + tY) = k(X + tY)^{-T}$
	$\frac{d}{dt} \log \det(X + tY) = \text{tr}((X + tY)^{-1} Y)$
	$\frac{d^2}{dt^2} \log \det(X + tY) = -\text{tr}((X + tY)^{-1} Y (X + tY)^{-1} Y)$
	$\frac{d}{dt} \log \det(X + tY)^{-1} = -\text{tr}((X + tY)^{-1} Y)$
	$\frac{d^2}{dt^2} \log \det(X + tY)^{-1} = \text{tr}((X + tY)^{-1} Y (X + tY)^{-1} Y)$
	$\begin{aligned} \frac{d}{dt} \log \det(\delta(A(x + tY) + a)^2 + \mu I) \\ = \text{tr}((\delta(A(x + tY) + a)^2 + \mu I)^{-1} 2\delta(A(x + tY) + a)\delta(Ay)) \end{aligned}$

### D.2.5 determinant

$$\begin{aligned}
 \nabla_X \det X &= \nabla_X \det X^T = \det(X)X^{-T} \\
 \nabla_X \det X^{-1} &= -\det(X^{-1})X^{-T} = -\det(X)^{-1}X^{-T} \\
 \nabla_X \det^\mu X &= \mu \det^\mu(X)X^{-T} \\
 \nabla_X \det X^\mu &= \mu \det(X^\mu)X^{-T} \\
 \nabla_X \det X^k &= k \det^{k-1}(X)(\text{tr}(X)I - X^T), \quad X \in \mathbb{R}^{2 \times 2} \\
 \nabla_X \det X^k &= \nabla_X \det^k X = k \det(X^k)X^{-T} = k \det^k(X)X^{-T} \\
 \nabla_X \det^\mu(X + tY) &= \mu \det^\mu(X + tY)(X + tY)^{-T} \\
 \nabla_X \det(X + tY)^k &= \nabla_X \det^k(X + tY) = k \det^k(X + tY)(X + tY)^{-T} \\
 \frac{d}{dt} \det(X + tY) &= \det(X + tY) \text{tr}((X + tY)^{-1}Y) \\
 \frac{d^2}{dt^2} \det(X + tY) &= \det(X + tY)(\text{tr}^2((X + tY)^{-1}Y) - \text{tr}((X + tY)^{-1}Y(X + tY)^{-1}Y)) \\
 \frac{d}{dt} \det(X + tY)^{-1} &= -\det(X + tY)^{-1} \text{tr}((X + tY)^{-1}Y) \\
 \frac{d^2}{dt^2} \det(X + tY)^{-1} &= \det(X + tY)^{-1}(\text{tr}^2((X + tY)^{-1}Y) + \text{tr}((X + tY)^{-1}Y(X + tY)^{-1}Y)) \\
 \frac{d}{dt} \det^\mu(X + tY) &= \mu \det^\mu(X + tY) \text{tr}((X + tY)^{-1}Y)
 \end{aligned}$$

### D.2.6 logarithmic

Matrix logarithm.

$$\begin{aligned}
 \frac{d}{dt} \log(X + tY)^\mu &= \mu Y(X + tY)^{-1} = \mu(X + tY)^{-1}Y, \quad XY = YX \\
 \frac{d}{dt} \log(I - tY)^\mu &= -\mu Y(I - tY)^{-1} = -\mu(I - tY)^{-1}Y \quad [229, \text{ p.493}]
 \end{aligned}$$

### D.2.7 exponential

Matrix exponential. [84, §3.6, §4.5] [368, §5.4]

$\nabla_X e^{\text{tr}(Y^T X)} = \nabla_X \det e^{Y^T X} = e^{\text{tr}(Y^T X)} Y$	$(\forall X, Y)$
$\nabla_X \text{tr } e^{YX} = e^{Y^T X^T} Y^T = Y^T e^{X^T Y^T}$	$(\forall X, Y)$
$\nabla_X \text{tr}(Ae^{YX}) = \dots$	
$\nabla_x \mathbf{1}^T e^{Ax} = A^T e^{Ax}$	
$\nabla_x \mathbf{1}^T e^{ Ax } = A^T \delta(\text{sgn}(Ax)) e^{ Ax }$	$(Ax)_i \neq 0$
$\nabla_x \log(\mathbf{1}^T e^x) = \frac{1}{\mathbf{1}^T e^x} e^x$	
$\nabla_x^2 \log(\mathbf{1}^T e^x) = \frac{1}{\mathbf{1}^T e^x} \left( \delta(e^x) - \frac{1}{\mathbf{1}^T e^x} e^x e^{x^T} \right)$	
$\nabla_x \prod_{i=1}^k x_i^{\frac{1}{k}} = \frac{1}{k} \left( \prod_{i=1}^k x_i^{\frac{1}{k}} \right) \mathbf{1}/x$	
$\nabla_x^2 \prod_{i=1}^k x_i^{\frac{1}{k}} = -\frac{1}{k} \left( \prod_{i=1}^k x_i^{\frac{1}{k}} \right) \left( \delta(x)^{-2} - \frac{1}{k} (\mathbf{1}/x)(\mathbf{1}/x)^T \right)$	
$\frac{d}{dt} e^{tY} = e^{tY} Y = Y e^{tY}$	
$\frac{d}{dt} e^{X+tY} = e^{X+tY} Y = Y e^{X+tY},$	$XY = YX$
$\frac{d^2}{dt^2} e^{X+tY} = e^{X+tY} Y^2 = Y e^{X+tY} Y = Y^2 e^{X+tY},$	$XY = YX$
$\frac{d^j}{dt^j} e^{\text{tr}(X+tY)} = e^{\text{tr}(X+tY)} \text{tr}^j(Y)$	

#### D.2.7.0.1 Exercise. Expand these tables.

Provide four unfinished table entries indicated by ... in §D.2.1 & §D.2.3. ▼

#### D.2.7.0.2 Exercise. $\log$ .

(§D.1.7, §3.5.4)

Find the first four terms of the Taylor series expansion for  $\log x$  about  $x=1$ . Plot the supporting hyperplane to the hypograph of  $\log x$  at  $\begin{bmatrix} x \\ \log x \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ . Prove  $\log x \leq x-1$ . ▼



## Appendix E

# Projection

*Rob Reiner's "A Few Good Men" is one of those movies that tells you what it's going to do, does it, and then tells you what it did. It doesn't think the audience is very bright.*

—Roger Ebert, 1992

For any  $t > 0$  [182, §2.1]

$$I - A(A^T A + t I)^{-1} A^T = t(AA^T + t I)^{-1} \quad (2045)$$

For any  $A \in \mathbb{R}^{m \times n}$ , the *pseudoinverse* [228, §7.3 prob.9] [280, §6.12 prob.19] [181, §5.5.4] [368, App.A]

$$A^\dagger \triangleq \lim_{t \rightarrow 0^+} (A^T A + t I)^{-1} A^T = \lim_{t \rightarrow 0^+} A^T (AA^T + t I)^{-1} \in \mathbb{R}^{n \times m} \quad (2046)$$

is a unique matrix from the convex optimal solution set to minimize<sub>X</sub>  $\|AX - I\|_F^2$  (§3.6.0.0.2). Pseudoinverse  $A^\dagger$  is that unique matrix satisfying the *Moore-Penrose conditions*: [230, §1.3] [436]

$$\begin{array}{lll} 1. & AA^\dagger A = A & 3. & (AA^\dagger)^T = AA^\dagger \\ 2. & A^\dagger AA^\dagger = A^\dagger & 4. & (A^\dagger A)^T = A^\dagger A \end{array} \quad (2047)$$

which are necessary and sufficient to establish the pseudoinverse whose principal action is to injectively map  $\mathcal{R}(A)$  onto  $\mathcal{R}(A^T)$  (Figure 187). Conditions 1 and 3 are necessary and sufficient for  $AA^\dagger$  to be the orthogonal projector on  $\mathcal{R}(A)$ , while conditions 2 and 4 hold iff  $A^\dagger A$  is the orthogonal projector on  $\mathcal{R}(A^T)$ .

Range and nullspace of the pseudoinverse [298] [364, §III.1 exer.1]

$$\mathcal{R}(A^\dagger) = \mathcal{R}(A^T), \quad \mathcal{R}(A^{\dagger T}) = \mathcal{R}(A) \quad (2048)$$

$$\mathcal{N}(A^\dagger) = \mathcal{N}(A^T), \quad \mathcal{N}(A^{\dagger T}) = \mathcal{N}(A) \quad (2049)$$

can be derived by singular value decomposition (§A.6).

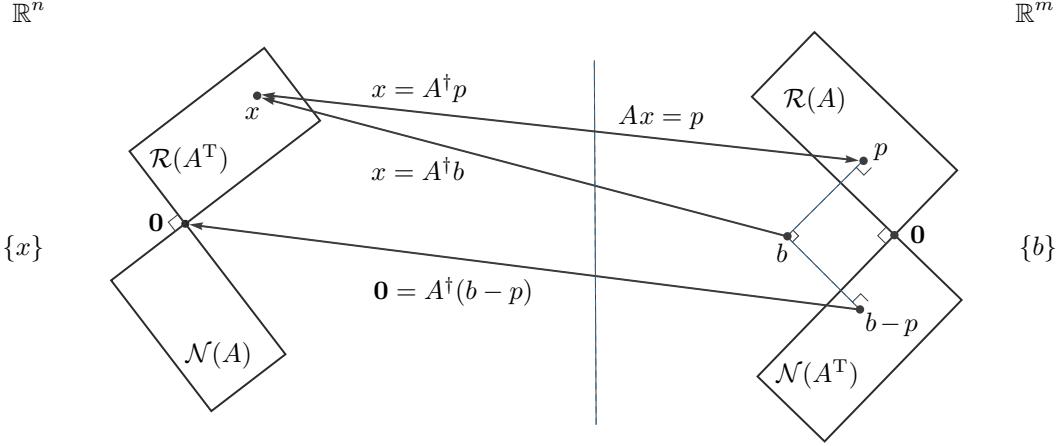


Figure 187: (confer Figure 18) Pseudoinverse  $A^\dagger \in \mathbb{R}^{n \times m}$  action: [368, p.449]  
 Component of  $b$  in  $\mathcal{N}(A^T)$  maps to  $\mathbf{0}$ , while component of  $b$  in  $\mathcal{R}(A)$  maps to rowspace  $\mathcal{R}(A^T)$ . For any  $A \in \mathbb{R}^{m \times n}$ ,  $p = AA^\dagger b$  and inversion is bijective  $\forall p \in \mathcal{R}(A)$ .  $x = A^\dagger b \Leftrightarrow x \in \mathcal{R}(A^T) \& b - Ax \perp \mathcal{R}(A) \Leftrightarrow x \perp \mathcal{N}(A) \& b - Ax \in \mathcal{N}(A^T)$ . [52]

The following relations reliably hold without qualification:

- a.  $A^{T\dagger} = A^{\dagger T}$
- b.  $A^{\dagger\dagger} = A$
- c.  $(AA^T)^\dagger = A^{\dagger T}A^\dagger$
- d.  $(A^TA)^\dagger = A^\dagger A^{\dagger T}$
- e.  $(AA^\dagger)^\dagger = AA^\dagger$
- f.  $(A^\dagger A)^\dagger = A^\dagger A$

Yet for arbitrary  $A, B$  it is generally true that  $(AB)^\dagger \neq B^\dagger A^\dagger$ :

**E.0.0.0.1 Theorem.** *Pseudoinverse of product.* [193] [61] [270, exer.7.23]  
 For  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times k}$

$$(AB)^\dagger = B^\dagger A^\dagger \quad (2050)$$

if and only if

$$\mathcal{R}(A^T AB) \subseteq \mathcal{R}(B) \quad \text{and} \quad \mathcal{R}(BB^T A^T) \subseteq \mathcal{R}(A^T) \quad (2051)$$

◇

Pseudoinverse of normalized vector  $u$  is the vector transpose. Otherwise,

$$u^\dagger = \frac{u^T}{\|u\|^2} \quad (2052)$$

$U^\dagger = U^T$  for orthonormal (including the orthogonal) matrices  $U$ . So, for orthonormal matrices  $U, Q$  and arbitrary  $A$

$$(UAQ^T)^\dagger = QA^\dagger U^T \quad (2053)$$

**E.0.0.0.2 Exercise.** *Kronecker pseudoinverse.*

Prove:

$$(A \otimes B)^\dagger = A^\dagger \otimes B^\dagger \quad (2054)$$

▼

### E.0.1 Logical deductions

When  $A$  is invertible,  $A^\dagger = A^{-1}$ ; so  $I^\dagger = I$  and  $A^\dagger A = AA^\dagger = I$ .

More generally, for  $A \in \mathbb{R}^{m \times n}$  [174, §5.3.3.1] [270, §7] [330]

g. $A^\dagger A = I$ ,	$A^\dagger = (A^T A)^{-1} A^T$ ,	$\text{rank } A = n$
h. $AA^\dagger = I$ ,	$A^\dagger = A^T (A A^T)^{-1}$ ,	$\text{rank } A = m$
i. $A^\dagger A \omega = \omega$ ,		$\omega \in \mathcal{R}(A^T)$
j. $AA^\dagger v = v$ ,		$v \in \mathcal{R}(A)$
k. $A^\dagger A = AA^\dagger$ ,		$A$ normal
l. $A^{k\dagger} = A^{\dagger k}$ ,		$k$ an integer, $A$ normal

Equivalent to the corresponding Moore-Penrose condition (2047):

$$\begin{array}{lll} 1. & A^T = A^T A A^\dagger & \text{or} \\ 2. & A^{\dagger T} = A^{\dagger T} A^\dagger A & \text{or} \end{array} \quad \begin{array}{ll} A^T = A^\dagger A A^T & \\ A^{\dagger T} = A A^\dagger A^{\dagger T} & \end{array}$$

There exists a singularity in the definition of pseudoinverse:

$$\mathbf{0}^\dagger = \mathbf{0}^T, \quad \delta(\mathbf{0})^\dagger = \delta(\mathbf{0})^T \quad (2055)$$

0 on the main diagonal of a diagonal matrix remains 0 in the pseudoinverse, whereas nonzero entries are inverted. When  $A$  is symmetric,  $A^\dagger$  is symmetric and (§A.6)

$$A \succeq 0 \Leftrightarrow A^\dagger \succeq 0 \quad (2056)$$

For any  $A \in \mathbb{R}^{m \times n}$  (§E.3.1)

$$A^\dagger A \succeq 0, \quad AA^\dagger \succeq 0, \quad I - A^\dagger A \succeq 0, \quad I - AA^\dagger \succeq 0 \quad (2057)$$

#### E.0.1.0.1 Example. Solution to classical linear equation $Ax = b$ .

In §2.5.1.1, the solution set to matrix equation  $Ax = b$  was represented as an intersection of hyperplanes. Regardless of rank of  $A$  or its shape {thin, square, wide}, interpretation as a hyperplane intersection (describing a possibly empty affine set) generally holds. Unique solution occurs when the hyperplanes intersect at a single point; e.g, when  $A$  is invertible.

#### vector $b$ not in range of matrix $A$

Given arbitrary matrix  $A$  (of any rank and dimension) and vector  $b$  not necessarily in  $\mathcal{R}(A)$ , we wish to find a best solution  $x^*$  to

$$Ax \approx b \quad (2058)$$

in a Euclidean sense by solving an algebraic expression for orthogonal projection of  $b$  on  $\mathcal{R}(A)$

$$\underset{x}{\text{minimize}} \quad \|Ax - b\|_2^2 \quad (2059)$$

Necessary and sufficient condition for optimal solution to this unconstrained optimization is the so-called *normal equation* that results from zeroing the convex objective's gradient: (§D.2.1)

$$A^T A x = A^T b \quad (2060)$$

*normal* because error vector  $b - Ax$  is perpendicular to  $\mathcal{R}(A)$ ; *id est*,  $A^T(b - Ax) = \mathbf{0}$ . Given any matrix  $A$  and any vector  $b$ , the normal equation is solvable exactly; always so, because  $\mathcal{R}(A^T A) = \mathcal{R}(A^T)$  and  $A^T b \in \mathcal{R}(A^T)$ . Given particular  $x_p \in \mathcal{R}(A^T)$  solving

(2060), then it is necessarily unique in  $\mathcal{R}(A^T)$  (Figure 187) and  $x_p = x^* = A^\dagger b$ . When  $A$  is thin-or-square full-rank, normal equation (2060) can be solved exactly by inversion:

$$x^* = (A^T A)^{-1} A^T b \equiv A^\dagger b \quad (2061)$$

For matrix  $A$  of arbitrary rank and shape, on the other hand,  $A^T A$  might not be invertible. Yet the normal equation can always be solved exactly by: (2046)

$$x^* = \lim_{t \rightarrow 0^+} (A^T A + t I)^{-1} A^T b = A^\dagger b \quad (2062)$$

invertible for any positive value of  $t$  by (1605). The exact inversion (2061) and this pseudoinverse solution (2062) each solve the same limited regularization

$$\lim_{t \rightarrow 0^+} \underset{x}{\text{minimize}} \|Ax - b\|_2^2 + t \|x\|_2^2 \equiv \lim_{t \rightarrow 0^+} \underset{x}{\text{minimize}} \left\| \begin{bmatrix} A \\ \sqrt{t} I \end{bmatrix} x - \begin{bmatrix} b \\ \mathbf{0} \end{bmatrix} \right\|_2^2 \quad (2063)$$

simultaneously providing least squares solution to (2059) and the classical *least norm* solution<sup>E.1</sup> [368, App.A.4] [52]

$$\begin{aligned} & \underset{x}{\text{minimize}} \|x\|_2^2 \\ & \text{subject to } Ax = AA^\dagger b \end{aligned} \quad (2064)$$

where  $AA^\dagger b$  is the orthogonal projection of vector  $b$  on  $\mathcal{R}(A)$ . Least norm solution can be interpreted as orthogonal projection of the origin  $\mathbf{0}$  on affine subset  $\mathcal{A} = \{x \mid Ax = AA^\dagger b\}$ ; (§E.5.0.0.6, §E.5.0.0.7)

$$\begin{aligned} & \underset{x}{\text{minimize}} \|x - \mathbf{0}\|_2^2 \\ & \text{subject to } x \in \mathcal{A} \end{aligned} \quad (2065)$$

equivalently, maximization of the Euclidean ball until it kisses  $\mathcal{A}$ ; rather,  $\arg \text{dist}(\mathbf{0}, \mathcal{A})$ .

#### vector $b$ in range of matrix $A$

If matrix  $A$  is rank deficient or wide, then there exists an infinite number of exact solutions  $x$  when  $b \in \mathcal{R}(A)$  and many ways to find them: for arbitrary choice of norm here, (§3.2.0.0.1)

$$\begin{aligned} & \text{find } x \\ & \text{subject to } Ax = b \end{aligned} \quad (1)$$

$$\underset{x}{\text{minimize}} \|Ax - b\| \quad (2)$$

$$\underset{x}{\text{minimize}} x^T A^T A x - 2x^T A^T b \quad (3) \quad (2066)$$

$$\begin{aligned} & \text{find } x \\ & \text{subject to } A^T A x = A^T b \end{aligned} \quad (4)$$

$$\underset{x}{\text{minimize}} \|A^T A x - A^T b\| \quad (5)$$

Depending upon numerical method of solution, each formulation can produce a different result. But each solution can always be decomposed into a sum of two vectors: a unique vector  $x_p$  from rowspace  $\mathcal{R}(A^T)$ , the other vector  $\eta$  from nullspace  $\mathcal{N}(A)$ :

$$x^* = x_p + \eta \quad (2067)$$

where  $x_p = A^\dagger b$  is the solution of least Euclidean norm:

$$\begin{aligned} & \underset{x}{\text{minimize}} \|x\|_2^2 \\ & \text{subject to } Ax = b \end{aligned} \quad (2068)$$

□

---

<sup>E.1</sup>This means: optimal solutions of lesser norm [*sic*] than the so-called *least norm* solution (2064) can be obtained (at expense of approximation  $Ax \approx b$ ; hence, of perpendicularity) by ignoring the limiting operation and introducing finite positive values of  $t$  into (2063).

## E.1 Idempotent matrices

Projection matrices are square and defined by *idempotence*,  $P^2 = P$ ; [368, §2.6] [230, §1.3] equivalent to the condition,  $P$  be diagonalizable [228, §3.3 prob.3] with eigenvalues  $\phi_i \in \{0, 1\}$ . [455, §4.1 thm.4.1] Idempotent matrices are not necessarily symmetric. The transpose of an idempotent matrix remains idempotent;  $P^T P^T = P^T$ . Solely excepting  $P = I$ , all projection matrices are neither orthogonal (§B.5) or invertible. [368, §3.4] The collection of all projection matrices of particular dimension does not form a convex set.

Suppose we wish to project nonorthogonally (*obliquely*) on the range of any particular matrix  $A \in \mathbb{R}^{m \times n}$ . All idempotent matrices projecting nonorthogonally on  $\mathcal{R}(A)$  may be expressed: (confer(2093))

$$P = A(A^\dagger + BZ^T) \in \mathbb{R}^{m \times m} \quad (2069)$$

where  $\mathcal{R}(P) = \mathcal{R}(A)$ , <sup>E.2</sup>  $B \in \mathbb{R}^{n \times k}$  for positive integer  $k$  is arbitrary, and  $Z \in \mathbb{R}^{m \times k}$  is any matrix whose range is in  $\mathcal{N}(A^T)$ ; *id est*,

$$A^T Z = A^\dagger Z = \mathbf{0} \quad (2070)$$

Evidently, the collection of nonorthogonal projectors projecting on  $\mathcal{R}(A)$  is an affine subset

$$\mathcal{P}_k = \left\{ A(A^\dagger + BZ^T) \mid B \in \mathbb{R}^{n \times k} \right\} \quad (2071)$$

When matrix  $A$  in (2069) is thin full-rank ( $A^\dagger A = I$ ) or has orthonormal columns ( $A^T A = I$ ), either property leads to a biorthogonal characterization of nonorthogonal projection:

### E.1.1 Biorthogonal characterization of projector

Any nonorthogonal projector  $P^2 = P \in \mathbb{R}^{m \times m}$  projecting on nontrivial  $\mathcal{R}(U)$  can be defined by a biorthogonality condition  $Q^T U = I$ ; the *biorthogonal decomposition* of  $P$  being (confer(2069))<sup>E.3</sup>

$$P = U Q^T, \quad Q^T U = I \quad (2072)$$

where<sup>E.4</sup> (§B.1.1.1)

$$\mathcal{R}(P) = \mathcal{R}(U), \quad \mathcal{N}(P) = \mathcal{N}(Q^T) \quad (2073)$$

and where generally (confer(2100))<sup>E.5</sup>

$$P^2 = P, \quad P^T \neq P, \quad P^\dagger \neq P, \quad \|P\|_2 \neq 1, \quad P \not\succeq 0 \quad (2074)$$

and  $P$  is not nonexpansive (2101) (2292).  $P^2 = P$  is necessary and sufficient.

<sup>E.2</sup>**Proof.**  $\mathcal{R}(P) \subseteq \mathcal{R}(A)$  is obvious [368, §3.6]. By (143) and (144),

$$\begin{aligned} \mathcal{R}(A^\dagger + BZ^T) &= \{(A^\dagger + BZ^T)y \mid y \in \mathbb{R}^m\} \\ &\supseteq \{(A^\dagger + BZ^T)y \mid y \in \mathcal{R}(A)\} = \mathcal{R}(A^T) \end{aligned}$$

$$\begin{aligned} \mathcal{R}(P) &= \{A(A^\dagger + BZ^T)y \mid y \in \mathbb{R}^m\} \\ &\supseteq \{A(A^\dagger + BZ^T)y \mid (A^\dagger + BZ^T)y \in \mathcal{R}(A^T)\} = \mathcal{R}(A) \end{aligned} \quad \blacklozenge$$

<sup>E.3</sup>  $A \leftarrow U$ ,  $A^\dagger + BZ^T \leftarrow Q^T$

<sup>E.4</sup>**Proof.** Obviously,  $\mathcal{R}(P) \subseteq \mathcal{R}(U)$ . Because  $Q^T U = I$

$$\begin{aligned} \mathcal{R}(P) &= \{UQ^T x \mid x \in \mathbb{R}^m\} \\ &\supseteq \{UQ^T Uy \mid y \in \mathbb{R}^k\} = \mathcal{R}(U) \end{aligned} \quad \blacklozenge$$

<sup>E.5</sup> Orthonormal decomposition (2097) (confer §E.3.4) is a special case of biorthogonal decomposition (2072) characterized by (2100). So, these characteristics (2074) are not necessary conditions for biorthogonality.

( $\Leftarrow$ ) To verify assertion (2072) we observe: because idempotent matrices are diagonalizable (§A.5), [228, §3.3 prob.3] they must have the form (1699)

$$P = S\Phi S^{-1} = \sum_{i=1}^m \phi_i s_i w_i^T = \sum_{i=1}^{k \leq m} s_i w_i^T \quad (2075)$$

that is a sum of  $k = \text{rank } P$  independent projector dyads (idempotent dyads, §B.1.1, §E.6.2.1) where  $\phi_i \in \{0, 1\}$  are the eigenvalues of  $P$  [455, §4.1 thm.4.1] in diagonal matrix  $\Phi \in \mathbb{R}^{m \times m}$  arranged in nonincreasing order, and where  $s_i, w_i \in \mathbb{R}^m$  are the right- and left-eigenvectors of  $P$ , respectively, which are independent and real.<sup>E.6</sup> Therefore

$$U \triangleq S(:, 1:k) = [s_1 \cdots s_k] \in \mathbb{R}^{m \times k} \quad (2076)$$

is the full-rank matrix  $S \in \mathbb{R}^{m \times m}$  having  $m - k$  columns (corresponding to 0 eigenvalues) truncated, while

$$Q^T \triangleq S^{-1}(1:k, :) = \begin{bmatrix} w_1^T \\ \vdots \\ w_k^T \end{bmatrix} \in \mathbb{R}^{k \times m} \quad (2077)$$

is matrix  $S^{-1}$  having the corresponding  $m - k$  rows truncated. By the 0 eigenvalues theorem (§A.7.3.0.1),  $\mathcal{R}(U) = \mathcal{R}(P)$ ,  $\mathcal{R}(Q) = \mathcal{R}(P^T)$ , and

$$\begin{aligned} \mathcal{R}(P) &= \text{span}\{s_i \mid \phi_i = 1 \ \forall i\} \\ \mathcal{N}(P) &= \text{span}\{s_i \mid \phi_i = 0 \ \forall i\} \\ \mathcal{R}(P^T) &= \text{span}\{w_i \mid \phi_i = 1 \ \forall i\} \\ \mathcal{N}(P^T) &= \text{span}\{w_i \mid \phi_i = 0 \ \forall i\} \end{aligned} \quad (2078)$$

Thus biorthogonality  $Q^T U = I$  is a necessary condition for idempotence, and so the collection of nonorthogonal projectors projecting on  $\mathcal{R}(U)$  is the affine subset  $\mathcal{P}_k = U \mathcal{Q}_k^T$  where  $\mathcal{Q}_k = \{Q \mid Q^T U = I, Q \in \mathbb{R}^{m \times k}\}$ .

( $\Rightarrow$ ) Biorthogonality is a sufficient condition for idempotence;

$$P^2 = \sum_{i=1}^k s_i w_i^T \sum_{j=1}^k s_j w_j^T = P \quad (2079)$$

*id est*, if the cross-products are annihilated, then  $P^2 = P$ .  $\blacklozenge$

Nonorthogonal projection of  $x$  on  $\mathcal{R}(P)$  has expression like a biorthogonal expansion,

$$Px = U Q^T x = \sum_{i=1}^k w_i^T x s_i \quad (2080)$$

When the domain is restricted to range of  $P$ , say  $x = U\xi$  for  $\xi \in \mathbb{R}^k$ , then  $x = Px = U Q^T U \xi = U\xi$  and expansion is unique because the eigenvectors are linearly independent. Otherwise, any component of  $x$  in  $\mathcal{N}(P) = \mathcal{N}(Q^T)$  will be annihilated. Direction of nonorthogonal projection is orthogonal to  $\mathcal{R}(Q) \Leftrightarrow Q^T U = I$ ; *id est*, for  $Px \in \mathcal{R}(U)$  (confer (2090))

$$Px - x \perp \mathcal{R}(Q) \text{ in } \mathbb{R}^m \quad (2081)$$

---

<sup>E.6</sup>Eigenvectors of a real matrix corresponding to real eigenvalues must be real. (§A.5.0.0.1)

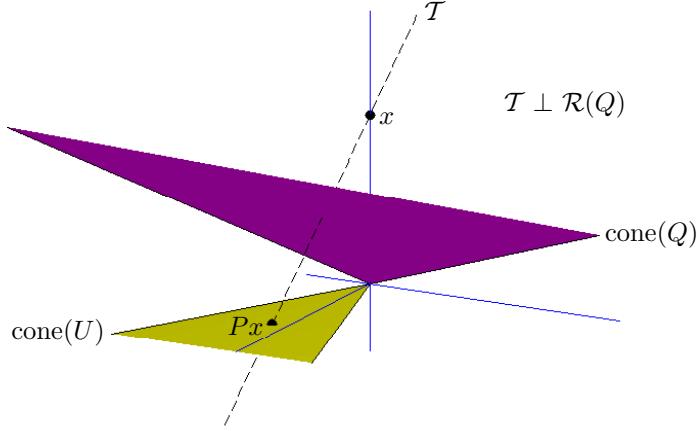


Figure 188: Nonorthogonal projection of  $x \in \mathbb{R}^3$  on  $\mathcal{R}(U) = \mathbb{R}^2$  under biorthogonality condition; *id est*,  $Px = UQ^T x$  such that  $Q^T U = I$ . Any point along imaginary line  $T$  connecting  $x$  to  $Px$  will be projected nonorthogonally on  $Px$  with respect to horizontal plane constituting  $\mathbb{R}^2 = \text{aff } \text{cone}(U)$  in this example. Extreme directions of  $\text{cone}(U)$  correspond to two columns of  $U$ ; likewise for  $\text{cone}(Q)$ . For purpose of illustration, we truncate each conic hull by truncating coefficients of conic combination at unity. Conic hull  $\text{cone}(Q)$  is headed upward at an angle, out of plane of page. Nonorthogonal projection would fail were  $\mathcal{N}(Q^T)$  in  $\mathcal{R}(U)$  (were  $T$  parallel to a line in  $\mathcal{R}(U)$ ).

#### E.1.1.0.1 Example. Illustration of nonorthogonal projector.

Figure 188 shows  $\text{cone}(U)$ , conic hull of the columns of

$$U = \begin{bmatrix} 1 & 1 \\ -0.5 & 0.3 \\ 0 & 0 \end{bmatrix} \quad (2082)$$

from nonorthogonal projector  $P = UQ^T$ . Matrix  $U$  has a limitless number of left inverses because  $\mathcal{N}(U^T)$  is nontrivial. Similarly depicted is left inverse  $Q^T$  from (2069)

$$\begin{aligned} Q &= U^{\dagger T} + ZB^T = \begin{bmatrix} 0.3750 & 0.6250 \\ -1.2500 & 1.2500 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} [0.5 \ 0.5] \\ &= \begin{bmatrix} 0.3750 & 0.6250 \\ -1.2500 & 1.2500 \\ 0.5000 & 0.5000 \end{bmatrix} \end{aligned} \quad (2083)$$

where  $Z \in \mathcal{N}(U^T)$  and matrix  $B$  is selected arbitrarily; *id est*,  $Q^T U = I$  because  $U$  is full-rank.

Direction of projection on  $\mathcal{R}(U)$  is orthogonal to  $\mathcal{R}(Q)$ . Any point along line  $T$  in the figure, for example, will have the same projection. Were matrix  $Z$  instead equal to  $\mathbf{0}$ , then  $\text{cone}(Q)$  would become the relative dual to  $\text{cone}(U)$  (sharing the same affine hull; §2.13.9, confer Figure 60a). In that case, projection  $Px = UU^\dagger x$  of  $x$  on  $\mathcal{R}(U)$  becomes orthogonal projection (and unique minimum-distance).  $\square$

### E.1.2 Idempotence summary

Nonorthogonal subspace-projector  $P$  is a (convex) linear operator defined by idempotence or biorthogonal decomposition (2072), but characterized not by symmetry or positive semidefiniteness or nonexpansivity (2101).

## E.2 $I - P$ , Projection on algebraic complement

It follows from the diagonalizability of idempotent matrices that  $I - P$  must also be a projection matrix because it too is idempotent, and because it may be expressed

$$I - P = S(I - \Phi)S^{-1} = \sum_{i=1}^m (1 - \phi_i) s_i w_i^T \quad (2084)$$

where  $(1 - \phi_i) \in \{1, 0\}$  are the eigenvalues of  $I - P$  (1606) whose eigenvectors  $s_i, w_i$  are identical to those of  $P$  in (2075). A consequence of that complementary relationship of eigenvalues is the fact, [381, §2] [375, §2] for subspace projector  $P = P^2 \in \mathbb{R}^{m \times m}$

$$\begin{aligned} \mathcal{R}(P) &= \text{span}\{s_i \mid \phi_i = 1 \ \forall i\} = \text{span}\{s_i \mid (1 - \phi_i) = 0 \ \forall i\} = \mathcal{N}(I - P) \\ \mathcal{N}(P) &= \text{span}\{s_i \mid \phi_i = 0 \ \forall i\} = \text{span}\{s_i \mid (1 - \phi_i) = 1 \ \forall i\} = \mathcal{R}(I - P) \\ \mathcal{R}(P^T) &= \text{span}\{w_i \mid \phi_i = 1 \ \forall i\} = \text{span}\{w_i \mid (1 - \phi_i) = 0 \ \forall i\} = \mathcal{N}(I - P^T) \\ \mathcal{N}(P^T) &= \text{span}\{w_i \mid \phi_i = 0 \ \forall i\} = \text{span}\{w_i \mid (1 - \phi_i) = 1 \ \forall i\} = \mathcal{R}(I - P^T) \end{aligned} \quad (2085)$$

that is easy to see from (2075) and (2084). Idempotent  $I - P$  therefore projects vectors on its range:  $\mathcal{N}(P)$ . Because all eigenvectors of a real idempotent matrix are real and independent, the algebraic complement of  $\mathcal{R}(P)$  [254, §3.3] is equivalent to  $\mathcal{N}(P)$ ; [E.7](#)  
*id est,*

$$\mathcal{R}(P) \oplus \mathcal{N}(P) = \mathcal{R}(P^T) \oplus \mathcal{N}(P^T) = \mathcal{R}(P^T) \oplus \mathcal{N}(P) = \mathcal{R}(P) \oplus \mathcal{N}(P^T) = \mathbb{R}^m \quad (2086)$$

because  $\mathcal{R}(P) \oplus \mathcal{R}(I - P) = \mathbb{R}^m$ . For idempotent  $P \in \mathbb{R}^{m \times m}$ , consequently,

$$\text{rank } P + \text{rank}(I - P) = m \quad (2087)$$

**E.2.0.0.1 Theorem.** *Rank Trace.* [455, §4.1 prob.9] (confer §E.3.2.0.1)

$$\begin{aligned} P^2 &= P \\ \Leftrightarrow \\ \text{rank } P &= \text{tr } P \quad \text{and} \quad \text{rank}(I - P) = \text{tr}(I - P) \end{aligned} \quad (2088)$$

◇

### E.2.1 Universal projector characteristic

Although projection is not necessarily orthogonal and  $\mathcal{R}(P) \not\perp \mathcal{R}(I - P)$  in general, still for any projector  $P$  and any  $x \in \mathbb{R}^m$

$$Px + (I - P)x = x \quad (2089)$$

must hold where  $\mathcal{R}(I - P) = \mathcal{N}(P)$  is the algebraic complement of  $\mathcal{R}(P)$ . Algebraic complement of closed convex cone  $\mathcal{K}$ , for example, is the negative dual cone  $-\mathcal{K}^*$ . (2224) (Figure 193)

---

[E.7](#)Same phenomenon occurs with symmetric (nonidempotent) matrices, for example. When summands in  $A \oplus B = \mathbb{R}^m$  are orthogonal vector spaces, the algebraic complement is the orthogonal complement.

### E.3 Symmetric idempotent matrices

When idempotent matrix  $P$  is symmetric,  $P$  is an orthogonal projector. In other words, the direction of projection of point  $x \in \mathbb{R}^m$  on subspace  $\mathcal{R}(P)$  is orthogonal to  $\mathcal{R}(P)$ ; *id est*, for  $P^2 = P \in \mathbb{S}^m$  and projection  $Px \in \mathcal{R}(P)$

$$Px - x \perp \mathcal{R}(P) \text{ in } \mathbb{R}^m \quad (2090)$$

(confer(2081)) Perpendicularity is a necessary and sufficient condition for orthogonal projection on a subspace. [122, §4.9]

A condition equivalent to (2090) is: Norm of direction  $x - Px$  is the infimum over all nonorthogonal projections of  $x$  on  $\mathcal{R}(P)$ ; [280, §3.3] for  $P^2 = P \in \mathbb{S}^m$ ,  $\mathcal{R}(P) = \mathcal{R}(A)$ , matrices  $A, B, Z$  and positive integer  $k$  as defined for (2069), and given  $x \in \mathbb{R}^m$

$$\|x - Px\|_2 = \|x - AA^\dagger x\|_2 = \inf_{B \in \mathbb{R}^{n \times k}} \|x - A(A^\dagger + BZ^T)x\|_2 = \text{dist}(x, \mathcal{R}(P)) \quad (2091)$$

The infimum is attained for  $\mathcal{R}(B) \subseteq \mathcal{N}(A)$  over any affine subset of nonorthogonal projectors (2071) indexed by  $k$ .

Proof is straightforward: The vector 2-norm is a convex function. Setting gradient of the norm-square to  $\mathbf{0}$ , applying §D.2,

$$\begin{aligned} (A^T A B Z^T - A^T (I - A A^\dagger)) x x^T A &= \mathbf{0} \\ \Leftrightarrow \\ A^T A B Z^T x x^T A &= \mathbf{0} \end{aligned} \quad (2092)$$

because  $A^T = A^T A A^\dagger$ . Projector  $P = A A^\dagger$  is therefore unique; the minimum-distance projector is the orthogonal projector, and *vice versa*. ♦

We get (confer(2069))

$$P = A A^\dagger \in \mathbb{R}^{m \times m} \quad (2093)$$

so this projection matrix must be symmetric. Then for any matrix  $A \in \mathbb{R}^{m \times n}$ , symmetric idempotent  $P$  projects a given vector  $x$  in  $\mathbb{R}^m$  orthogonally on  $\mathcal{R}(A)$ . Vector  $A^\dagger x$  in  $\mathbb{R}^n$  comprises coefficients of projection. Under either condition (2090) or (2091), the projection  $Px$  is unique minimum-distance; for subspaces, perpendicularity and minimum-distance conditions are equivalent.

#### E.3.1 Four subspaces

We summarize the orthogonal projectors projecting on the four fundamental subspaces: for  $A \in \mathbb{R}^{m \times n}$

$$\begin{array}{lll} A^\dagger A : \mathbb{R}^n & \text{on} & \mathcal{R}(A^\dagger A) = \mathcal{R}(A^T) = \mathcal{R}(A^T A) \\ A A^\dagger : \mathbb{R}^m & \text{on} & \mathcal{R}(A A^\dagger) = \mathcal{R}(A) = \mathcal{R}(A A^T) \\ I - A^\dagger A : \mathbb{R}^n & \text{on} & \mathcal{R}(I - A^\dagger A) = \mathcal{N}(A) = \mathcal{R}(I - A^T A) \\ I - A A^\dagger : \mathbb{R}^m & \text{on} & \mathcal{R}(I - A A^\dagger) = \mathcal{N}(A^T) = \mathcal{R}(I - A A^T) \end{array} \quad (2094)$$

Given a known subspace, matrix  $A$  is neither unique or necessarily full-rank. Despite that, a basis for each fundamental subspace comprises the linearly independent column vectors from its associated symmetric projection matrix:

$$\begin{array}{lll} \text{basis } \mathcal{R}(A^T) & \subseteq A^\dagger A & \subseteq \mathcal{R}(A^T) \\ \text{basis } \mathcal{R}(A) & \subseteq A A^\dagger & \subseteq \mathcal{R}(A) \\ \text{basis } \mathcal{N}(A) & \subseteq I - A^\dagger A & \subseteq \mathcal{N}(A) \\ \text{basis } \mathcal{N}(A^T) & \subseteq I - A A^\dagger & \subseteq \mathcal{N}(A^T) \end{array} \quad (2095)$$

For completeness:[E.8](#) ([2085](#))

$$\begin{aligned}\mathcal{N}(A^\dagger A) &= \mathcal{N}(A) = \mathcal{N}(A^T A) \\ \mathcal{N}(AA^\dagger) &= \mathcal{N}(A^T) = \mathcal{N}(AA^T) \\ \mathcal{N}(I - A^\dagger A) &= \mathcal{R}(A^T) = \mathcal{N}(I - A^T A) \\ \mathcal{N}(I - AA^\dagger) &= \mathcal{R}(A) = \mathcal{N}(I - AA^T)\end{aligned}\tag{2096}$$

### E.3.2 Orthogonal characterization of projector

Any symmetric projector  $P^2 = P \in \mathbb{S}^m$ , projecting on nontrivial  $\mathcal{R}(Q)$ , can be defined by orthonormality condition  $Q^T Q = I$ . When thin matrix  $Q \in \mathbb{R}^{m \times k}$  is orthonormal (has orthonormal columns), then  $Q^\dagger = Q^T$  by the Moore-Penrose conditions ([2047](#)). Hence, any  $P$  having an *orthonormal decomposition* ([§E.3.4](#))

$$P = QQ^T, \quad Q^T Q = I \tag{2097}$$

where [[368](#), §3.3] ([1765](#))

$$\mathcal{R}(P) = \mathcal{R}(Q), \quad \mathcal{N}(P) = \mathcal{N}(Q^T) \tag{2098}$$

is an orthogonal projector projecting on  $\mathcal{R}(Q)$ ; for  $Px \in \mathcal{R}(Q)$  (*confer* ([2081](#)))

$$Px - x \perp \mathcal{R}(Q) \text{ in } \mathbb{R}^m \tag{2099}$$

From ([2097](#)), orthogonal projector  $P$  is obviously positive semidefinite ([§A.3.1.0.6](#)); necessarily, (*confer* ([2074](#)))

$$P^2 = P, \quad P^T = P, \quad P^\dagger = P, \quad \|P\|_2 = 1, \quad P \succeq 0 \tag{2100}$$

and  $\|Px\| = \|QQ^T x\| = \|Q^T x\|$  because  $\|Qy\| = \|y\| \forall y \in \mathbb{R}^k$ .  $P^2 = P$  is also sufficient. All orthogonal projectors are therefore *nonexpansive* because

$$\sqrt{\langle Px, x \rangle} = \|Px\| = \|Q^T x\| \leq \|x\| \quad \forall x \in \mathbb{R}^m \tag{2101}$$

the Bessel inequality, [[122](#)] [[254](#)] with equality when  $x \in \mathcal{R}(Q)$ .

From the diagonalization of idempotent matrices ([2075](#)) on page [572](#)

$$P = S\Phi S^T = \sum_{i=1}^m \phi_i s_i s_i^T = \sum_{i=1}^{k \leq m} s_i s_i^T \tag{2102}$$

orthogonal projection of point  $x$  on  $\mathcal{R}(P)$  has expression like an orthogonal expansion [[122](#), §4.10]

$$Px = QQ^T x = \sum_{i=1}^k s_i^T x s_i \tag{2103}$$

where

$$Q = S(:, 1:k) = [s_1 \cdots s_k] \in \mathbb{R}^{m \times k} \tag{2104}$$

and where the  $s_i$  [*sic*] are orthonormal eigenvectors, of symmetric idempotent  $P$ , corresponding to eigenvalues  $\phi_i \in \{0, 1\}$ . When the domain is restricted to range of  $P$ , say  $x = Q\xi$  for  $\xi \in \mathbb{R}^k$ , then  $x = Px = QQ^T Q\xi = Q\xi$  and expansion is unique. Otherwise, any component of  $x$  in  $\mathcal{N}(Q^T)$  will be annihilated.

---

**E.8 Proof** is by singular value decomposition ([§A.6.2](#)):  $\mathcal{N}(A^\dagger A) \subseteq \mathcal{N}(A)$  is obvious. Conversely, suppose  $A^\dagger A x = \mathbf{0}$ . Then  $x^T A^\dagger A x = x^T Q Q^T x = \|Q^T x\|^2 = 0$  where  $A = U\Sigma Q^T$  is the subcompact SVD. Because  $\mathcal{R}(Q) = \mathcal{R}(A^T)$ , then  $x \in \mathcal{N}(A)$  which implies  $\mathcal{N}(A^\dagger A) \supseteq \mathcal{N}(A)$ .  $\therefore \mathcal{N}(A^\dagger A) = \mathcal{N}(A)$ . ♦

**E.3.2.0.1 Theorem.** *Symmetric rank trace.* (confer §E.2.0.0.1, (1610))

$$\begin{aligned} P^T &= P, \quad P^2 = P \\ &\Leftrightarrow \\ \text{rank } P &= \text{tr } P = \|P\|_F^2 \quad \text{and} \quad \text{rank}(I - P) = \text{tr}(I - P) = \|I - P\|_F^2 \end{aligned} \quad (2105)$$

◊

**Proof.** We take, as given, Theorem E.2.0.0.1 establishing idempotence. We have left only to show  $\text{tr } P = \|P\|_F^2 \Rightarrow P^T = P$ , established in [455, §7.1]. ◆

### E.3.3 Summary, symmetric idempotent

(confer §E.1.2) Orthogonal projector  $P$  is a (convex) linear operator defined [225, §A.3.1] by idempotence and symmetry, and characterized by positive semidefiniteness and nonexpansivity. The algebraic complement (§E.2) to  $\mathcal{R}(P)$  becomes the orthogonal complement  $\mathcal{R}(I - P)$ ; *id est*,  $\mathcal{R}(P) \perp \mathcal{R}(I - P)$ .

### E.3.4 Orthonormal decomposition

When  $Z = \mathbf{0}$  in the general nonorthogonal projector  $A(A^\dagger + BZ^T)$  (2069), an orthogonal projector results (for any matrix  $A$ ) characterized principally by idempotence and symmetry. Any real orthogonal projector may, in fact, be represented by an orthonormal decomposition such as (2097). [230, §1 prob.42]

To verify that assertion for the four fundamental subspaces (2094), we need only to express  $A$  by subcompact singular value decomposition (§A.6.2): From pseudoinverse (1736) of  $A = U\Sigma Q^T \in \mathbb{R}^{m \times n}$

$$\begin{aligned} AA^\dagger &= U\Sigma\Sigma^\dagger U^T = UU^T, & A^\dagger A &= Q\Sigma^\dagger\Sigma Q^T = QQ^T \\ I - AA^\dagger &= I - UU^T = U^\perp U^{\perp T}, & I - A^\dagger A &= I - QQ^T = Q^\perp Q^{\perp T} \end{aligned} \quad (2106)$$

where  $U^\perp \in \mathbb{R}^{m \times m - \text{rank } A}$  holds columnar an orthonormal basis for the orthogonal complement of  $\mathcal{R}(U)$ , and likewise for  $Q^\perp \in \mathbb{R}^{n \times n - \text{rank } A}$ . Existence of an orthonormal decomposition is sufficient to establish idempotence and symmetry of an orthogonal projector (2097). ◆

### E.3.5 Unifying trait of all projectors: direction

Whereas nonorthogonal projectors possess only a biorthogonal decomposition (§E.1.1), relation (2106) shows: orthogonal projectors simultaneously possess a biorthogonal decomposition ( $AA^\dagger$  whence  $Px = AA^\dagger x$ ) and an orthonormal decomposition ( $UU^T$  whence  $Px = UU^T x$ ). Orthogonal projection of a point is unique but its expansion is not; *e.g.*,  $A$  can have dependent columns.

#### E.3.5.1 orthogonal projector, orthonormal decomposition

Consider orthogonal expansion of  $x \in \mathcal{R}(U)$ :

$$x = UU^T x = \sum_{i=1}^n u_i u_i^T x \quad (2107)$$

a sum of one-dimensional orthogonal projections (§E.6.3) where

$$U \triangleq [u_1 \cdots u_n] \quad \text{and} \quad U^T U = I \quad (2108)$$

and where the subspace projector has two expressions: (2106)

$$AA^\dagger \triangleq UU^T \quad (2109)$$

where  $A \in \mathbb{R}^{m \times n}$  has rank  $n$ . The direction of projection of  $x$  on  $u_j$  for some  $j \in \{1 \dots n\}$ , for example, is orthogonal to  $u_j$  but parallel to a vector in the span of all remaining vectors constituting the columns of  $U$ ;

$$\begin{aligned} u_j^T(u_j u_j^T x - x) &= 0 \\ u_j u_j^T x - x &= u_j u_j^T x - UU^T x \in \mathcal{R}(\{u_i \mid i=1 \dots n, i \neq j\}) \end{aligned} \quad (2110)$$

### E.3.5.2 orthogonal projector, biorthogonal decomposition

We get a similar result for biorthogonal expansion of  $x \in \mathcal{R}(A)$ . Define

$$A \triangleq [a_1 \ a_2 \ \cdots \ a_n] \in \mathbb{R}^{m \times n} \quad (2111)$$

and rows of the pseudoinverse<sup>E.9</sup>

$$A^\dagger \triangleq \begin{bmatrix} a_1^{*\top} \\ a_2^{*\top} \\ \vdots \\ a_n^{*\top} \end{bmatrix} \in \mathbb{R}^{n \times m} \quad (2112)$$

under biorthogonality condition  $A^\dagger A = I$ . In biorthogonal expansion (§2.13.9)

$$x = AA^\dagger x = \sum_{i=1}^n a_i a_i^{*\top} x \quad (2113)$$

the direction of projection of  $x$  on  $a_j$  for some particular  $j \in \{1 \dots n\}$ , for example, is orthogonal to  $a_j^{*\top}$  and parallel to a vector in the span of all the remaining vectors constituting the columns of  $A$ ;

$$\begin{aligned} a_j^{*\top}(a_j a_j^{*\top} x - x) &= 0 \\ a_j a_j^{*\top} x - x &= a_j a_j^{*\top} x - AA^\dagger x \in \mathcal{R}(\{a_i \mid i=1 \dots n, i \neq j\}) \end{aligned} \quad (2114)$$

### E.3.5.3 nonorthogonal projector, biorthogonal decomposition

Because the result in §E.3.5.2 is independent of matrix symmetry  $AA^\dagger = (AA^\dagger)^T$ , we must get the same result for any nonorthogonal projector characterized by a biorthogonality condition; namely, for nonorthogonal projector  $P = UQ^T$  (2072) under biorthogonality condition  $Q^T U = I$ , in biorthogonal expansion of  $x \in \mathcal{R}(U)$

$$x = UQ^T x = \sum_{i=1}^k u_i q_i^T x \quad (2115)$$

where

$$\begin{aligned} U &\triangleq [u_1 \cdots u_k] \in \mathbb{R}^{m \times k} \\ Q^T &\triangleq \begin{bmatrix} q_1^T \\ \vdots \\ q_k^T \end{bmatrix} \in \mathbb{R}^{k \times m} \end{aligned} \quad (2116)$$

---

<sup>E.9</sup>Notation \* in this context connotes extreme direction of a dual cone; e.g., (413) or Example E.5.0.0.3.

the direction of projection of  $x$  on  $u_j$  is orthogonal to  $q_j$  and parallel to a vector in the span of the remaining  $u_i$ :

$$\begin{aligned} q_j^T(u_j q_j^T x - x) &= 0 \\ u_j q_j^T x - x &= u_j q_j^T x - UQ^T x \in \mathcal{R}(\{u_i \mid i=1 \dots k, i \neq j\}) \end{aligned} \quad (2117)$$

## E.4 Algebra of projection on affine subsets

Let  $P_{\mathcal{A}}x$  denote orthogonal or nonorthogonal projection of  $x$  on affine subset  $\mathcal{A} \triangleq \mathcal{R} + \alpha$  where  $\mathcal{R}$  is a subspace and  $\alpha \in \mathcal{A}$ . Let  $P_{\mathcal{R}}x$  denote projection of  $x$  on  $\mathcal{R}$  in the same direction. Then, because  $\mathcal{R}$  is parallel to  $\mathcal{A}$ , it holds:

$$\begin{aligned} P_{\mathcal{A}}x &= P_{\mathcal{R}+\alpha}x = P_{\mathcal{R}}x + (I - P_{\mathcal{R}})\alpha \\ &= P_{\mathcal{R}}(x - \alpha) + \alpha \end{aligned} \quad (2118)$$

Orthogonal or nonorthogonal subspace-projector  $P_{\mathcal{R}}$  is a linear operator ( $P_{\mathcal{A}}$  is not).  $P_{\mathcal{R}}(x+y)=P_{\mathcal{R}}x$  whenever  $y \perp \mathcal{R}$  and  $P_{\mathcal{R}}$  is an orthogonal projector.

**E.4.0.0.1 Theorem.** *Orthogonal projection on affine subset.* [122, §9.26]

Let  $\mathcal{A} = \mathcal{R} + \alpha$  be an affine subset where  $\alpha \in \mathcal{A}$ , and let  $\mathcal{R}^\perp$  be the orthogonal complement of subspace  $\mathcal{R}$ . Then  $P_{\mathcal{A}}x$  is the orthogonal projection of  $x$  on  $\mathcal{A}$  if and only if

$$P_{\mathcal{A}}x \in \mathcal{A}, \quad \langle P_{\mathcal{A}}x - x, a - \alpha \rangle = 0 \quad \forall a \in \mathcal{A} \quad (2119)$$

or if and only if

$$P_{\mathcal{A}}x \in \mathcal{A}, \quad P_{\mathcal{A}}x - x \in \mathcal{R}^\perp \quad (2120)$$

◇

**E.4.0.0.2 Example.** *Intersection of affine subsets.*

When designing an optimization problem, intersection of sets is easy to express when the individual sets themselves are. Suppose  $\mathcal{A} = \{x \mid Ax = b\}$  and  $\mathcal{C} = \{x \mid Cx = d\}$  denote two affine subsets. Then membership to their intersection

$$\mathcal{A} \cap \mathcal{C} = \left\{ x \mid \begin{bmatrix} A \\ C \end{bmatrix} x = \begin{bmatrix} b \\ d \end{bmatrix} \right\} \quad (2121)$$

is realized as solution to simultaneous equations. In 1937, Kaczmarz instead proposed alternating projection (§E.10) on hyperplanes (whose normals constitute the rows of matrices  $A$  and  $C$ ) as a means to overcome numerical instability when solving very large systems. The intersection may then be described as fixed points of projection on affine subsets, assuming  $\mathcal{A} \cap \mathcal{C} \neq \emptyset$

$$\mathcal{A} \cap \mathcal{C} = \{x \mid P_{\mathcal{C}}P_{\mathcal{A}}x = x\} \quad (2122)$$

an affine system of equations, by (2118), extensible to  $M$  intersecting hyperplanes (§E.5.0.0.5):

$$\mathcal{A} \cap \mathcal{C} = \{x \mid P_M \cdots P_1 x = x\} \quad (2123)$$

□

## E.5 Projection examples

### E.5.0.0.1 Example. Orthogonal projection on orthogonal basis.

Orthogonal projection on a subspace can instead be accomplished by orthogonally projecting on the individual members of an orthogonal basis for that subspace. Suppose, for example, matrix  $A \in \mathbb{R}^{m \times n}$  holds an orthonormal basis for  $\mathcal{R}(A)$  in its columns;  $A \triangleq [a_1 \ a_2 \ \cdots \ a_n]$ . Then orthogonal projection of vector  $x \in \mathbb{R}^n$  on  $\mathcal{R}(A)$  is a sum of one-dimensional orthogonal projections

$$Px = AA^\dagger x = A(A^\top A)^{-1}A^\top x = AA^\top x = \sum_{i=1}^n a_i a_i^\top x \quad (2124)$$

where each symmetric dyad  $a_i a_i^\top$  is an orthogonal projector projecting on  $\mathcal{R}(a_i)$ . (§E.6.3) Because  $\|x - Px\|$  is minimized by orthogonal projection,  $Px$  is considered to be the best approximation (in the Euclidean sense) to  $x$  from the set  $\mathcal{R}(A)$ . [122, §4.9]  $\square$

### E.5.0.0.2 Example. Orthogonal projection on span of nonorthogonal basis.

Orthogonal projection on a subspace can also be accomplished by projecting nonorthogonally on the individual members of any nonorthogonal basis for that subspace. This interpretation is in fact the principal application of the pseudoinverse we discussed. Now suppose matrix  $A$  holds a nonorthogonal basis for  $\mathcal{R}(A)$  in its columns,

$$A = [a_1 \ a_2 \ \cdots \ a_n] \in \mathbb{R}^{m \times n} \quad (2111)$$

and define the rows  $a_i^{*\top}$  of its pseudoinverse  $A^\dagger$  as in (2112). Then orthogonal projection of vector  $x \in \mathbb{R}^n$  on  $\mathcal{R}(A)$  is a sum of one-dimensional nonorthogonal projections

$$Px = AA^\dagger x = \sum_{i=1}^n a_i a_i^{*\top} x \quad (2125)$$

where each nonsymmetric dyad  $a_i a_i^{*\top}$  is a nonorthogonal projector projecting on  $\mathcal{R}(a_i)$ , (§E.6.1) idempotent because of biorthogonality condition  $A^\dagger A = I$ .

The projection  $Px$  is regarded as the best approximation to  $x$  from the set  $\mathcal{R}(A)$ , as it was in Example E.5.0.0.1.  $\square$

### E.5.0.0.3 Example. Biorthogonal expansion as nonorthogonal projection.

Biorthogonal expansion can be viewed as a sum of components, each a nonorthogonal projection on the range of an extreme direction of a pointed polyhedral cone  $\mathcal{K}$ ; e.g, Figure 189.

Suppose matrix  $A \in \mathbb{R}^{m \times n}$  holds a nonorthogonal basis for  $\mathcal{R}(A)$  in its columns as in (2111), and the rows of pseudoinverse  $A^\dagger$  are defined as in (2112). Assuming the most general biorthogonality condition  $(A^\dagger + BZ^\top)A = I$  with  $BZ^\top$  defined as for (2069), then biorthogonal expansion of vector  $x$  is a sum of one-dimensional nonorthogonal projections; for  $x \in \mathcal{R}(A)$

$$x = A(A^\dagger + BZ^\top)x = AA^\dagger x = \sum_{i=1}^n a_i a_i^{*\top} x \quad (2126)$$

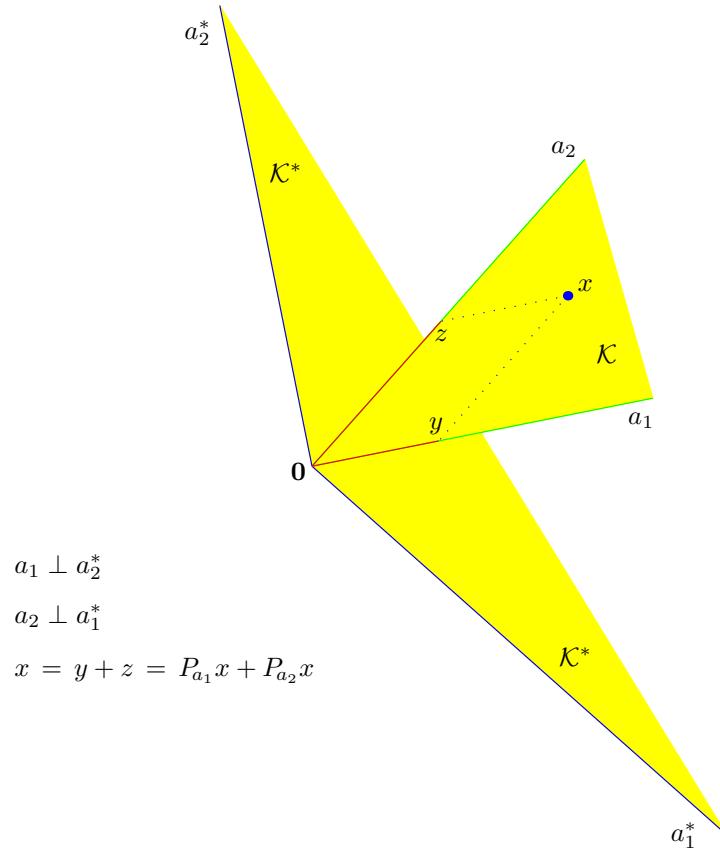


Figure 189: (confer Figure 67) Biorthogonal expansion of point  $x \in \text{aff } \mathcal{K}$  is found by projecting  $x$  nonorthogonally on extreme directions of polyhedral cone  $\mathcal{K} \subset \mathbb{R}^2$ . (Dotted lines of projection bound this translated negated cone.) Direction of projection  $P_{a_1}$  on extreme direction  $a_1$  is orthogonal to extreme direction  $a_1^*$  of dual cone  $\mathcal{K}^*$  and parallel to  $a_2$  (§E.3.5); similarly, direction of projection  $P_{a_2}$  on  $a_2$  is orthogonal to  $a_2^*$  and parallel to  $a_1$ . Point  $x$  is sum of nonorthogonal projections:  $x$  on  $\mathcal{R}(a_1)$  and  $x$  on  $\mathcal{R}(a_2)$ . Expansion is unique because extreme directions of  $\mathcal{K}$  are linearly independent. Were  $a_1$  orthogonal to  $a_2$ , then  $\mathcal{K}$  would be identical to  $\mathcal{K}^*$  and nonorthogonal projections would become orthogonal.

where each dyad  $a_i a_i^{*T}$  is a nonorthogonal projector projecting on  $\mathcal{R}(a_i)$ . (§E.6.1) The extreme directions of  $\mathcal{K} = \text{cone}(A)$  are  $\{a_1 \dots a_n\}$  the linearly independent columns of  $A$ , while the extreme directions  $\{a_1^* \dots a_n^*\}$  of relative dual cone  $\mathcal{K}^* \cap \text{aff } \mathcal{K} = \text{cone}(A^{\dagger T})$  (§2.13.10.4) correspond to the linearly independent (§B.1.1.1) rows of  $A^\dagger$ . Directions of nonorthogonal projection are determined by the pseudoinverse; *id est*, direction of projection  $a_i a_i^{*T}x - x$  on  $\mathcal{R}(a_i)$  is orthogonal to  $a_i^*$ .<sup>E.10</sup>

Because the extreme directions of this cone  $\mathcal{K}$  are linearly independent, component projections are unique in the sense:

- There is only one linear combination of extreme directions of  $\mathcal{K}$  that yields a particular point  $x \in \mathcal{R}(A)$  whenever

$$\mathcal{R}(A) = \text{aff } \mathcal{K} = \mathcal{R}(a_1) \oplus \mathcal{R}(a_2) \oplus \dots \oplus \mathcal{R}(a_n) \quad (2127)$$

□

#### E.5.0.0.4 Example. Nonorthogonal projection on elementary matrix.

Suppose  $P_Y$  is a linear nonorthogonal projector projecting on subspace  $\mathcal{Y}$ , and suppose the range of a vector  $u$  is linearly independent of  $\mathcal{Y}$ ; *id est*, for some other subspace  $\mathcal{M}$  containing  $\mathcal{Y}$  suppose

$$\mathcal{M} = \mathcal{R}(u) \oplus \mathcal{Y} \quad (2128)$$

Assuming  $P_M x = P_u x + P_Y x$  holds, then it follows for vector  $x \in \mathcal{M}$

$$P_u x = x - P_Y x, \quad P_Y x = x - P_u x \quad (2129)$$

nonorthogonal projection of  $x$  on  $\mathcal{R}(u)$  can be determined from nonorthogonal projection of  $x$  on  $\mathcal{Y}$ , and *vice versa*.

Such a scenario is realizable were there some arbitrary basis for  $\mathcal{Y}$  populating a full-rank thin-or-square matrix  $A$

$$A \triangleq [\text{basis } \mathcal{Y} \quad u] \in \mathbb{R}^{N \times n+1} \quad (2130)$$

With  $P_u = A(:, n+1)A^{\dagger}(n+1, :)$  and  $P_Y = A(:, 1:n)A^{\dagger}(1:n, :)$ , then  $P_M = AA^{\dagger}$  fulfills the requirements. Observe,  $P_M$  is an orthogonal projector whereas  $P_Y$  and  $P_u$  are nonorthogonal projectors.

Now suppose, for example,  $P_Y$  is an elementary matrix (§B.3); in particular,

$$P_Y = I - e_1 \mathbf{1}^T = \begin{bmatrix} \mathbf{0} & \sqrt{2}V_N \end{bmatrix} \in \mathbb{R}^{N \times N} \quad (2131)$$

where  $\mathcal{Y} = \mathcal{N}(\mathbf{1}^T)$ . We have  $\mathcal{M} = \mathbb{R}^N$ ,  $A = [\sqrt{2}V_N \quad e_1]$ , and  $u = e_1$ . Thus  $P_u = e_1 \mathbf{1}^T$  is a nonorthogonal projector projecting on  $\mathcal{R}(u)$  in a direction parallel to a vector in  $\mathcal{Y}$  (§E.3.5), and  $P_Y x = x - e_1 \mathbf{1}^T x$  is a nonorthogonal projection of  $x$  on  $\mathcal{Y}$  in a direction parallel to  $u$ . □

---

<sup>E.10</sup>This remains true in high dimension although only a little more difficult to visualize in  $\mathbb{R}^3$ ; *confer*, Figure 68.

**E.5.0.0.5 Example.** *Projecting on hyperplane, halfspace, slab.*

- Given hyperplane representation having  $b \in \mathbb{R}$  and nonzero normal  $a \in \mathbb{R}^n$

$$\partial\mathcal{H} = \{y \mid a^T y = b\} \subset \mathbb{R}^n \quad (116)$$

orthogonal projection of any point  $x \in \mathbb{R}^n$  on that hyperplane [233, §3.1] is unique:

$$\begin{aligned} Px &= x - a(a^T a)^{-1}(a^T x - b) \\ &= (I - aa^\dagger)x + aa^\dagger y_p \end{aligned} \quad (2132)$$

where  $y_p$  is any particular solution to  $a^T y = b$ .

- Orthogonal projection of  $x$  on the halfspace parametrized by  $b \in \mathbb{R}$  and nonzero normal  $a \in \mathbb{R}^n$

$$\mathcal{H}_- = \{y \mid a^T y \leq b\} \subset \mathbb{R}^n \quad (108)$$

is the point

$$Px = x - a(a^T a)^{-1} \max\{0, a^T x - b\} \quad (2133)$$

- Orthogonal projection of  $x$  on the convex slab (Figure 13), for  $c < b$

$$\{y \mid c \leq a^T y \leq b\} \subset \mathbb{R}^n \quad (2134)$$

is the point [168, §5.1]

$$Px = x - a(a^T a)^{-1} (\max\{0, a^T x - b\} - \max\{0, c - a^T x\}) \quad (2135)$$

□

**E.5.0.0.6 Example.** *Projecting origin on a hyperplane.* (confer §2.4.2.0.2)

Given the hyperplane representation having  $b \in \mathbb{R}$  and nonzero normal  $a \in \mathbb{R}^n$

$$\partial\mathcal{H} = \{y \mid a^T y = b\} \subset \mathbb{R}^n \quad (116)$$

orthogonal projection of the origin  $P\mathbf{0}$  on that hyperplane is the unique optimal solution to a minimization problem: (2091)

$$\begin{aligned} \|P\mathbf{0} - \mathbf{0}\|_2 &= \inf_{y \in \partial\mathcal{H}} \|y - \mathbf{0}\|_2 \\ &= \inf_{\xi \in \mathbb{R}^{n-1}} \|Z\xi + x\|_2 \end{aligned} \quad (2136)$$

where  $x$  is any solution to  $a^T y = b$ , and where the columns of  $Z \in \mathbb{R}^{n \times n-1}$  constitute a basis for  $\mathcal{N}(a^T)$  so that  $y = Z\xi + x \in \partial\mathcal{H}$  for all  $\xi \in \mathbb{R}^{n-1}$ .

The infimum can be found by setting the gradient (with respect to  $\xi$ ) of the strictly convex norm-square to  $\mathbf{0}$ . We find the minimizing argument

$$\xi^* = -(Z^T Z)^{-1} Z^T x \quad (2137)$$

so

$$y^* = (I - Z(Z^T Z)^{-1} Z^T)x \quad (2138)$$

and from (2094)

$$P\mathbf{0} = y^* = a(a^T a)^{-1} a^T x = \frac{a}{\|a\|} \frac{a^T}{\|a\|} x = \frac{a}{\|a\|^2} a^T x \triangleq A^\dagger A x = a \frac{b}{\|a\|^2} \quad (2139)$$

In words, any point  $x$  in the hyperplane  $\partial\mathcal{H}$  projected on its normal  $a$  (confer (2168)) yields that point  $y^*$  in the hyperplane closest to the origin. □

**E.5.0.0.7 Example.** *Projection on affine subset.*

The technique of Example E.5.0.0.6 is extensible. Given an intersection of hyperplanes

$$\mathcal{A} = \{y \mid Ay = b\} \subset \mathbb{R}^n \quad (151)$$

where each row of  $A \in \mathbb{R}^{m \times n}$  is nonzero and  $b \in \mathcal{R}(A)$ , then the orthogonal projection  $Px$  of any point  $x \in \mathbb{R}^n$  on  $\mathcal{A}$  is the solution to a minimization problem:

$$\begin{aligned} \|Px - x\|_2 &= \inf_{y \in \mathcal{A}} \|y - x\|_2 \\ &= \inf_{\xi \in \mathbb{R}^{n-\text{rank } A}} \|Z\xi + y_p - x\|_2 \end{aligned} \quad (2140)$$

where  $y_p$  is any solution to  $Ay = b$ , and where the columns of  $Z \in \mathbb{R}^{n \times n-\text{rank } A}$  constitute a basis for  $\mathcal{N}(A)$  so that  $y = Z\xi + y_p \in \mathcal{A}$  for all  $\xi \in \mathbb{R}^{n-\text{rank } A}$ . When rank of wide matrix  $A$  is  $n-1$ , then  $\mathcal{A}$  describes a line; when  $\text{rank } A = 1$ , then  $\mathcal{A}$  describes a hyperplane.

The infimum is found by setting the gradient of the strictly convex norm-square to  $\mathbf{0}$ . The minimizing argument is

$$\xi^* = -(Z^T Z)^{-1} Z^T (y_p - x) \quad (2141)$$

so

$$y^* = (I - Z(Z^T Z)^{-1} Z^T)(y_p - x) + x \quad (2142)$$

and from (2094),

$$\begin{aligned} Px = y^* &= x - A^\dagger(Ax - b) \\ &= (I - A^\dagger A)x + A^\dagger A y_p \end{aligned} \quad (2143)$$

which is a projection of  $x$  on  $\mathcal{N}(A)$  then translated perpendicularly with respect to the nullspace until it meets the affine subset  $\mathcal{A}$ .  $A^\dagger = A^T(AA^T)^{-1}$  when  $A$  is wide full-rank (§E.0.1).  $\square$

**E.5.0.0.8 Example.** *Projection on affine subset, vertex-description.*

Suppose now we instead describe the affine subset  $\mathcal{A}$  in terms of some given minimal set of generators arranged columnar in  $X \in \mathbb{R}^{n \times N}$  (77); *id est*,

$$\mathcal{A} = \text{aff } X = \{Xa \mid a^T \mathbf{1} = 1\} \subseteq \mathbb{R}^n \quad (79)$$

Here *minimal set* means  $XV_N = [x_2 - x_1 \ x_3 - x_1 \ \dots \ x_N - x_1]/\sqrt{2}$  (1105) is full-rank (§2.4.2.2) where  $V_N \in \mathbb{R}^{N \times N-1}$  is the Schoenberg auxiliary matrix (§B.4.2). For  $N=2$  affinely independent generators,  $\mathcal{A}$  represents a line; for  $N=n$  affinely independent generators,  $\mathcal{A}$  represents a hyperplane. Then the orthogonal projection  $Px$  of any point  $x \in \mathbb{R}^n$  on  $\mathcal{A}$  is the solution to a minimization problem:

$$\begin{aligned} \|Px - x\|_2 &= \inf_{a^T \mathbf{1} = 1} \|Xa - x\|_2 \\ &= \inf_{\xi \in \mathbb{R}^{N-1}} \|X(V_N \xi + a_p) - x\|_2 \end{aligned} \quad (2144)$$

where  $a_p$  is any solution to  $a^T \mathbf{1} = 1$ . We find the minimizing argument

$$\xi^* = -(V_N^T X^T X V_N)^{-1} V_N^T X^T (Xa_p - x) \quad (2145)$$

and so the orthogonal projection is [232, §3]

$$Px = Xa^* = (I - X V_N (X V_N)^\dagger) X a_p + X V_N (X V_N)^\dagger x \quad (2146)$$

a projection of point  $x$  on  $\mathcal{R}(X V_N)$  then translated perpendicularly with respect to that range until it meets the affine subset  $\mathcal{A}$ .  $\square$

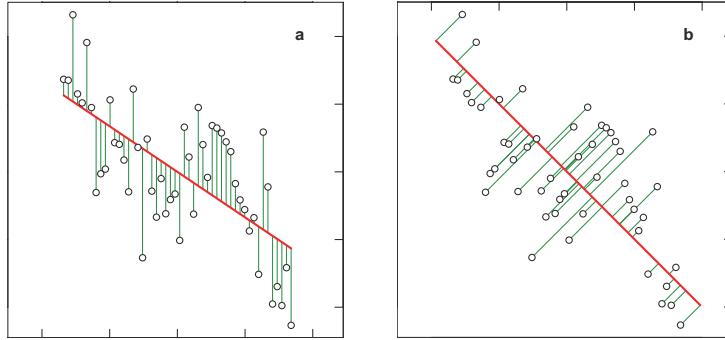


Figure 190: Line segment represents best affine fit to same set in  $\mathbb{R}^2$  by (a) *linear regression*, (b) *principal component analysis* (PCA, §E.5.0.0.9). Linear regression resembles PCA but differs insofar as regression minimizes distances of points in range of a real affine function (“vertical” distance). (Drawing by [Gavin Simpson](#).)

#### E.5.0.0.9 Example. Line fit to point cloud.

Principal component analysis (PCA) has found application to *machine learning*. One interpretation of PCA is to find a line maximizing Euclidean distances between (variance of) projections of given points on that line. A dual interpretation of PCA (Figure 190b) finds that line best fitting the same set of points by minimizing distances from points to line. Traditionally, PCA is carried out under the 2-norm; minimum Euclidean distance projection. But Brooks, Dulá, and Boone show that projection under 1-norm provides a better distance measure when outliers are problematic. [70]

In two dimensions, a line is a hyperplane. Orthogonal projection  $Px$  of point  $x$  on a hyperplane is described in Example E.5.0.0.5. Euclidean distance from hyperplane

$$\partial\mathcal{H} = \{y \mid a^T y = b\} \subset \mathbb{R}^n \quad (116)$$

to a point  $x \in \mathbb{R}^n$  is

$$\text{dist}(Px, x) = \frac{1}{\|a\|_2} |a^T x - b| \quad (2147)$$

Given a *cloud* of  $N$  points  $X = [x_1 \cdots x_N] \in \mathbb{R}^{n \times N}$ , minimization of distances is stated

$$\begin{aligned} &\underset{a \in \mathbb{R}^n, b \in \mathbb{R}}{\text{minimize}} && |a^T X - b \mathbf{1}^T| \mathbf{1} \\ &\text{subject to} && \|a\|_2 = 1 \end{aligned} \quad (2148)$$

having nonconvex constraint but solvable by convex iteration as in Example 4.7.0.0.1. Although projection on the hyperplane is orthogonal under 2-norm here, the objective is a 1-norm having lower bound equal to zero which occurs in the circumstance that all points in the cloud belong to the same hyperplane.

In higher dimension, a line is an intersection of  $n-1$  independent hyperplanes

$$\mathcal{A} = \{y \mid Ay = b\} \subset \mathbb{R}^n \quad (151)$$

where  $A \in \mathbb{R}^{n-1 \times n}$  is full-rank. Euclidean distance from  $\mathcal{A}$  to a given point  $x$  is (§E.5.0.0.7)

$$\text{dist}(P_{\mathcal{A}}x, x) = \|A^\dagger A(x - y_p)\|_2 \quad (2149)$$

where  $y_p$  is any point on the line. This says distance between line  $\mathcal{A}$  and point  $x$  is the same as projecting difference  $x - y_p$  on hyperplane  $\mathcal{R}(A^T)$  and then measuring length of the projected difference vector.

But minimizing distances to a point cloud that way is difficult. Instead we observe: any line in  $\mathbb{R}^n$  is orthogonal to a hyperplane

$$\mathcal{R} \triangleq \{y \mid a^T y = 0\} \subset \mathbb{R}^n \quad (2150)$$

containing the origin; an  $n-1$ -dimensional subspace of  $\mathbb{R}^n$ . Then the line is described

$$\mathcal{A} = \mathcal{R}^\perp + y_p = \{a\xi + y_p \mid \xi \in \mathbb{R}\} \quad (2151)$$

From algebra of projection (2118)

$$P_{\mathcal{A}}x = P_{\mathcal{R}^\perp + y_p}x = P_{\mathcal{R}^\perp}(x - y_p) + y_p \quad (2152)$$

and from projection on orthogonal complement (2223)

$$P_{\mathcal{A}}x = x - P_{\mathcal{R}}(x - y_p) = x - P_{\mathcal{R}}x + y_p \quad (2153)$$

where  $y_p$  is assumed to be the point of line/subspace intersection.

$$\text{dist}(P_{\mathcal{A}}x, x) = \text{dist}(x - P_{\mathcal{R}}(x - y_p), x) = \|P_{\mathcal{R}}(x - y_p)\|_2 = \|(I - aa^\dagger)(x - y_p)\|_2 \quad (2154)$$

where  $a^\dagger y_p = 0$  by (2150). This says: distance between line  $\mathcal{A}$  and point  $x$  is the same as length of difference vector  $x - y_p$  as it appears in its projection on subspace  $\mathcal{R}$ . Absolute distances to the point cloud are minimized, for  $P \triangleq I - aa^\dagger$

$$\begin{array}{ll} \underset{P \in \mathbb{S}^n, y_p \in \mathbb{R}^n}{\text{minimize}} & \sum_{i=1}^N \|Px_i - y_p\|_2 \\ \text{subject to} & \text{tr } P = n-1 \\ & \text{tr}(I - P) = 1 \\ & \text{rank } P = n-1 \\ & \text{rank}(I - P) = 1 \end{array} \quad (2155)$$

which has nonconvex feasible set necessary and sufficient for  $P$  to be a rank  $n-1$  projection matrix by Theorem E.2.0.0.1. Matrix symmetry provides orthogonal projection. Equivalent semidefinite program (§3.5.3)

$$\begin{array}{ll} \underset{P \in \mathbb{S}^n, y_p \in \mathbb{R}^n, t \in \mathbb{R}^N}{\text{minimize}} & \mathbf{1}^T t + (\langle W, P \rangle + \langle I - W, I - P \rangle) \lambda \\ \text{subject to} & \begin{bmatrix} t_i I & Px_i - y_p \\ (Px_i - y_p)^T & t_i \end{bmatrix} \succeq 0, \quad i = 1 \dots N \\ & \text{tr } P = n-1 \\ & \text{tr}(I - P) = 1 \end{array} \quad (2156)$$

handles rank constraints by choosing direction matrix  $W$  via convex iteration (§4.5.1).  $\mathbf{1}^T t$  represents 1-norm although projection remains orthogonal.  $\lambda$  is a positive scalar determined via bisection so that  $\langle W, P \rangle + \langle I - W, I - P \rangle$  just vanishes.

Conventional PCA would set  $y_p^* = \frac{1}{N} X \mathbf{1}$  to geometric center of the given point cloud, then identify a best fitting line as parallel to the principal eigenvector of  $XX^T \in \mathbb{R}^{n \times n}$ . Although less computationally intensive than (2156), the conventional approach generally produces a different line. Linear regression is to be preferred, over either approach to PCA, when data derives from a function (not geometry).  $\square$

## E.6 Vectorization interpretation Projection on a matrix

### E.6.1 Nonorthogonal projection on a vector

Nonorthogonal projection of vector  $x$  on the range of vector  $y$  is accomplished using a normalized dyad  $P_0$  (§B.1); *videlicet*,

$$\frac{\langle z, x \rangle}{\langle z, y \rangle} y = \frac{z^T x}{z^T y} y = \frac{y z^T}{z^T y} x \triangleq P_0 x \quad (2157)$$

where  $\langle z, x \rangle / \langle z, y \rangle$  is the coefficient of projection on  $y$ . Because  $P_0^2 = P_0$  and  $\mathcal{R}(P_0) = \mathcal{R}(y)$ , rank-one matrix  $P_0$  is a nonorthogonal projector dyad projecting on  $\mathcal{R}(y)$ . Direction of nonorthogonal projection is orthogonal to  $z$ ; *id est*,

$$P_0 x - x \perp \mathcal{R}(P_0^T) \quad (2158)$$

### E.6.2 Nonorthogonal projection on vectorized matrix

Formula (2157) is extensible. Given  $X, Y, Z \in \mathbb{R}^{m \times n}$ , we have a one-dimensional nonorthogonal projection of  $X$  in isomorphic  $\mathbb{R}^{mn}$  on the range of vectorized  $Y$ : (§2.2)

$$\frac{\langle Z, X \rangle}{\langle Z, Y \rangle} Y, \quad \langle Z, Y \rangle \neq 0 \quad (2159)$$

where  $\langle Z, X \rangle / \langle Z, Y \rangle$  is the coefficient of projection. The inequality accounts for the fact: projection on  $\mathcal{R}(\text{vec } Y)$  is in a direction orthogonal to  $\text{vec } Z$ . Projection is one-dimensional because vectorized  $Y$  represents a point in  $\mathbb{R}^{mn}$ .

#### E.6.2.1 Nonorthogonal projection on dyad

Now suppose we have nonorthogonal projector dyad

$$P_0 = \frac{y z^T}{z^T y} \in \mathbb{R}^{m \times m} \quad (2160)$$

Analogous to (2157), for  $X \in \mathbb{R}^{m \times n}$

$$P_0 X P_0 = \frac{y z^T}{z^T y} X \frac{y z^T}{z^T y} = \frac{z^T X y}{(z^T y)^2} y z^T = \frac{\langle z y^T, X \rangle}{\langle z y^T, y z^T \rangle} y z^T \quad (2161)$$

is a one-dimensional nonorthogonal projection of matrix  $X$  on the range of vectorized dyad  $P_0$ ; from which it follows:

$$P_0 X P_0 = \frac{z^T X y}{z^T y} \frac{y z^T}{z^T y} = \left\langle \frac{z y^T}{z^T y}, X \right\rangle \frac{y z^T}{z^T y} = \langle P_0^T, X \rangle P_0 = \frac{\langle P_0^T, X \rangle}{\langle P_0^T, P_0 \rangle} P_0 \quad (2162)$$

Yet this relationship between matrix product and vector inner-product only holds for a dyad projector. When nonsymmetric projector  $P_0$  is rank-one as in (2160), therefore,

$$\mathcal{R}(\text{vec } P_0 X P_0) = \mathcal{R}(\text{vec } P_0) \text{ in } \mathbb{R}^{m^2} \quad (2163)$$

and

$$P_0 X P_0 - X \perp P_0^T \text{ in } \mathbb{R}^{m^2} \quad (2164)$$

**E.6.2.1.1 Example.** *Eigenvalues  $\lambda$  as coefficients of nonorthogonal projection.*  
 (confer §E.6.4.1.1) Any diagonalization (§A.5)

$$X = S\Lambda S^{-1} = \sum_{i=1}^m \lambda_i s_i w_i^T \in \mathbb{R}^{m \times m} \quad (1699)$$

may be expressed as a sum of one-dimensional nonorthogonal projections of  $X$ , each on the range of a vectorized eigenmatrix  $P_j \triangleq s_j w_j^T$ ;

$$\begin{aligned} X &= \sum_{i,j=1}^m \langle (Se_i e_j^T S^{-1})^T, X \rangle Se_i e_j^T S^{-1} \\ &= \sum_{j=1}^m \langle (s_j w_j^T)^T, X \rangle s_j w_j^T + \sum_{\substack{i,j=1 \\ j \neq i}}^m \langle (Se_i e_j^T S^{-1})^T, S\Lambda S^{-1} \rangle Se_i e_j^T S^{-1} \\ &= \sum_{j=1}^m \langle (s_j w_j^T)^T, X \rangle s_j w_j^T \\ &\triangleq \sum_{j=1}^m \langle P_j^T, X \rangle P_j = \sum_{j=1}^m s_j w_j^T X s_j w_j^T = \sum_{j=1}^m P_j X P_j \\ &= \sum_{j=1}^m \lambda_j s_j w_j^T \end{aligned} \quad (2165)$$

This biorthogonal expansion of matrix  $X$  is a sum of nonorthogonal projections because the term outside the projection coefficient  $\langle \cdot \rangle$  is not identical to the inside-term. (§E.6.4) The eigenvalues  $\lambda_j$  are coefficients of nonorthogonal projection of  $X$ , while the remaining  $M(M-1)/2$  coefficients (for  $i \neq j$ ) are zeroed by projection. When  $P_j$  is rank-one as in (2165),

$$\mathcal{R}(\text{vec } P_j X P_j) = \mathcal{R}(\text{vec } s_j w_j^T) = \mathcal{R}(\text{vec } P_j) \text{ in } \mathbb{R}^{m^2} \quad (2166)$$

and

$$P_j X P_j - X \perp P_j^T \text{ in } \mathbb{R}^{m^2} \quad (2167)$$

Were matrix  $X$  symmetric, then its eigenmatrices would also be. So, the one-dimensional projections would become orthogonal. (§E.6.4.1.1)  $\square$

### E.6.3 Orthogonal projection on a vector

The formula for orthogonal projection of vector  $x$  on the range of vector  $y$  (*one-dimensional projection*) is basic *analytic geometry*; [14, §3.3] [368, §3.2] [406, §2.2] [442, §1-8]

$$\frac{\langle y, x \rangle}{\langle y, y \rangle} y = \frac{y^T x}{y^T y} y = \frac{y y^T}{y^T y} x \triangleq P_1 x \quad (2168)$$

where  $\langle y, x \rangle / \langle y, y \rangle$  is the coefficient of projection on  $\mathcal{R}(y)$ . An equivalent description is: Vector  $P_1 x$  is the orthogonal projection of vector  $x$  on  $\mathcal{R}(P_1) = \mathcal{R}(y)$ . Rank-one matrix  $P_1$  is a projection matrix because  $P_1^2 = P_1$ . The direction of projection is orthogonal

$$P_1 x - x \perp \mathcal{R}(P_1) \quad (2169)$$

because  $P_1^T = P_1$ .

### E.6.4 Orthogonal projection on a vectorized matrix

From (2168), given instead  $X, Y \in \mathbb{R}^{m \times n}$ , we have the one-dimensional orthogonal projection of matrix  $X$  in isomorphic  $\mathbb{R}^{mn}$  on the range of vectorized  $Y$ : (§2.2)

$$\frac{\langle Y, X \rangle}{\langle Y, Y \rangle} Y \quad (2170)$$

where  $\langle Y, X \rangle / \langle Y, Y \rangle$  is the coefficient of projection. For orthogonal projection, the term outside the vector inner-products  $\langle \cdot \rangle$  must be identical to the terms inside in three places. Projection is one-dimensional because  $Y$  describes a point in vector space  $\mathbb{R}^{mn}$ .

#### E.6.4.1 Orthogonal projection on dyad

There is opportunity for insight when  $Y$  is a dyad  $yz^T$  (§B.1): Instead given  $X \in \mathbb{R}^{m \times n}$ ,  $y \in \mathbb{R}^m$ , and  $z \in \mathbb{R}^n$

$$\frac{\langle yz^T, X \rangle}{\langle yz^T, yz^T \rangle} yz^T = \frac{y^T X z}{y^T y z^T z} yz^T \quad (2171)$$

is the one-dimensional orthogonal projection of  $X$  in isomorphic  $\mathbb{R}^{mn}$  on the range of vectorized  $yz^T$ . To reveal obscured symmetric projection matrices  $P_1$  and  $P_2$  we rewrite (2171):

$$\frac{y^T X z}{y^T y z^T z} yz^T = \frac{yy^T}{y^T y} X \frac{zz^T}{z^T z} \triangleq P_1 X P_2 \quad (2172)$$

So for projector dyads, projection (2172) is the orthogonal projection in  $\mathbb{R}^{mn}$  if and only if projectors  $P_1$  and  $P_2$  are symmetric; [E.11](#) in other words,

- for one-dimensional orthogonal projection on the range of a vectorized dyad  $yz^T$ , the term outside the vector inner-products  $\langle \cdot \rangle$  in (2171) must be identical to the terms inside in three places.

When  $P_1$  and  $P_2$  are rank-one symmetric projectors as in (2172), (37)

$$\mathcal{R}(\text{vec } P_1 X P_2) = \mathcal{R}(\text{vec } yz^T) \text{ in } \mathbb{R}^{mn} \quad (2173)$$

and

$$P_1 X P_2 - X \perp yz^T \text{ in } \mathbb{R}^{mn} \quad (2174)$$

When  $y=z$  then  $P_1=P_2=P_2^T$  and

$$P_1 X P_1 = \langle P_1, X \rangle P_1 = \frac{\langle P_1, X \rangle}{\langle P_1, P_1 \rangle} P_1 \quad (2175)$$

meaning,  $P_1 X P_1$  is equivalent to one-dimensional orthogonal projection of matrix  $X$  on the range of vectorized projector dyad  $P_1$ . Yet this relationship between matrix product and vector inner-product does not hold for general symmetric projector matrices.

##### E.6.4.1.1 Example. Eigenvalues $\lambda$ as coefficients of orthogonal projection.

(confer §E.6.2.1.1) Let  $\mathcal{C}$  represent any convex subset of subspace  $\mathbb{S}^M$ , and let  $\mathcal{C}_1$  be any element of  $\mathcal{C}$ . Then  $\mathcal{C}_1$  can be expressed as the orthogonal expansion

$$\mathcal{C}_1 = \sum_{i=1}^M \sum_{\substack{j=1 \\ j \geq i}}^M \langle E_{ij}, \mathcal{C}_1 \rangle E_{ij} \in \mathcal{C} \subset \mathbb{S}^M \quad (2176)$$

---

[E.11](#) For diagonalizable  $X \in \mathbb{R}^{m \times m}$  (§A.5), its orthogonal projection (in isomorphic  $\mathbb{R}^{m^2}$ ) on the range of vectorized  $yz^T \in \mathbb{R}^{m \times m}$  becomes:

$$P_1 X P_2 = \sum_{i=1}^m \lambda_i P_1 s_i w_i^T P_2$$

When  $\mathcal{R}(P_1) = \mathcal{R}(w_j)$  and  $\mathcal{R}(P_2) = \mathcal{R}(s_j)$ , the  $j^{\text{th}}$  dyad term from the diagonalization is isolated but only, in general, to within a scale factor because neither set of left or right eigenvectors is necessarily orthonormal unless  $X$  is a normal matrix [455, §3.2]. Yet when  $\mathcal{R}(P_2) = \mathcal{R}(s_k)$ ,  $k \neq j \in \{1 \dots m\}$ , then  $P_1 X P_2 = \mathbf{0}$ .

where  $E_{ij} \in \mathbb{S}^M$  is a member of the standard orthonormal basis for  $\mathbb{S}^M$  (60). This expansion is a sum of one-dimensional orthogonal projections of  $\mathcal{C}_1$ ; each projection on the range of a vectorized standard basis matrix. Vector inner-product  $\langle E_{ij}, \mathcal{C}_1 \rangle$  is the coefficient of projection of  $\text{svec } \mathcal{C}_1$  on  $\mathcal{R}(\text{svec } E_{ij})$ .

When  $\mathcal{C}_1$  is any member of a convex set  $\mathcal{C}$  whose dimension is  $L$ , *Carathéodory's theorem* [126] [343] [225] [42] [43] guarantees that no more than  $L+1$  affinely independent members from  $\mathcal{C}$  are required to faithfully represent  $\mathcal{C}_1$  by their linear combination.<sup>E.12</sup>

Dimension of  $\mathbb{S}^M$  is  $L=M(M+1)/2$  in isometrically isomorphic  $\mathbb{R}^{M(M+1)/2}$ . Yet because any symmetric matrix can be diagonalized, (§A.5.1)  $\mathcal{C}_1 \in \mathbb{S}^M$  is a linear combination of its  $M$  eigenmatrices  $q_i q_i^T$  (§A.5.0.3) weighted by its eigenvalues  $\lambda_i$ ;

$$\mathcal{C}_1 = Q\Lambda Q^T = \sum_{i=1}^M \lambda_i q_i q_i^T \quad (2177)$$

where  $\Lambda \in \mathbb{S}^M$  is a diagonal matrix having  $\delta(\Lambda)_i = \lambda_i$ , and where  $Q = [q_1 \cdots q_M]$  is an orthogonal matrix in  $\mathbb{R}^{M \times M}$  containing corresponding eigenvectors.

To derive eigenvalue decomposition (2177) from expansion (2176),  $M$  standard basis matrices  $E_{ij}$  are rotated (§B.5) into alignment with the  $M$  eigenmatrices  $q_i q_i^T$  of  $\mathcal{C}_1$  by applying a *similarity transformation*; [368, §5.6]

$$\{QE_{ij}Q^T\} = \left\{ \begin{array}{ll} q_i q_i^T, & i=j=1 \dots M \\ \frac{1}{\sqrt{2}}(q_i q_j^T + q_j q_i^T), & 1 \leq i < j \leq M \end{array} \right\} \quad (2178)$$

which remains an orthonormal basis for  $\mathbb{S}^M$ . Then remarkably

$$\begin{aligned} \mathcal{C}_1 &= \sum_{\substack{i,j=1 \\ j \geq i}}^M \langle QE_{ij}Q^T, \mathcal{C}_1 \rangle QE_{ij}Q^T \\ &= \sum_{i=1}^M \langle q_i q_i^T, \mathcal{C}_1 \rangle q_i q_i^T + \sum_{\substack{i,j=1 \\ j > i}}^M \langle QE_{ij}Q^T, Q\Lambda Q^T \rangle QE_{ij}Q^T \\ &= \sum_{i=1}^M \langle q_i q_i^T, \mathcal{C}_1 \rangle q_i q_i^T \\ &\triangleq \sum_{i=1}^M \langle P_i, \mathcal{C}_1 \rangle P_i = \sum_{i=1}^M q_i q_i^T \mathcal{C}_1 q_i q_i^T = \sum_{i=1}^M P_i \mathcal{C}_1 P_i \\ &= \sum_{i=1}^M \lambda_i q_i q_i^T \end{aligned} \quad (2179)$$

this orthogonal expansion becomes the diagonalization; still a sum of one-dimensional orthogonal projections. The eigenvalues

$$\lambda_i = \langle q_i q_i^T, \mathcal{C}_1 \rangle \quad (2180)$$

are clearly coefficients of projection of  $\mathcal{C}_1$  on the range of each vectorized eigenmatrix. (*confer* §E.6.2.1.1) The remaining  $M(M-1)/2$  coefficients ( $i \neq j$ ) are zeroed by projection. When  $P_i$  is rank-one symmetric as in (2179),

$$\mathcal{R}(\text{svec } P_i \mathcal{C}_1 P_i) = \mathcal{R}(\text{svec } q_i q_i^T) = \mathcal{R}(\text{svec } P_i) \text{ in } \mathbb{R}^{M(M+1)/2} \quad (2181)$$

and

$$P_i \mathcal{C}_1 P_i - \mathcal{C}_1 \perp P_i \text{ in } \mathbb{R}^{M(M+1)/2} \quad (2182)$$

□

---

<sup>E.12</sup>Carathéodory's theorem guarantees existence of a biorthogonal expansion for any element in  $\text{aff } \mathcal{C}$  when  $\mathcal{C}$  is any pointed closed convex cone.

#### E.6.4.2 Positive semidefiniteness test as orthogonal projection

For any given  $X \in \mathbb{R}^{m \times m}$  the familiar quadratic construct  $y^T X y \geq 0$ , over broad domain, is a fundamental test for positive semidefiniteness. (§A.2) It is a fact that  $y^T X y$  is always proportional to a coefficient of orthogonal projection; letting  $z$  in formula (2171) become  $y \in \mathbb{R}^m$ , then  $P_2 = P_1 = yy^T / y^T y = yy^T / \|yy^T\|_2$  (confer (1768)) and formula (2172) becomes

$$\frac{\langle yy^T, X \rangle}{\langle yy^T, yy^T \rangle} yy^T = \frac{y^T X y}{y^T y} \frac{yy^T}{y^T y} = \frac{y^T y}{y^T y} X \frac{yy^T}{y^T y} \triangleq P_1 X P_1 \quad (2183)$$

Product  $P_1 X P_1$  is the one-dimensional orthogonal projection of  $X$  in isomorphic  $\mathbb{R}^{m^2}$  on the range of vectorized  $P_1$  because, by (2170) for  $\text{rank } P_1 = 1$  and  $P_1^2 = P_1 \in \mathbb{S}^m$  (confer (2162))

$$P_1 X P_1 = \frac{y^T X y}{y^T y} \frac{yy^T}{y^T y} = \left\langle \frac{yy^T}{y^T y}, X \right\rangle \frac{yy^T}{y^T y} = \langle P_1, X \rangle P_1 = \frac{\langle P_1, X \rangle}{\langle P_1, P_1 \rangle} P_1 \quad (2184)$$

The coefficient of orthogonal projection  $\langle P_1, X \rangle = y^T X y / (y^T y)$  is also known as *Rayleigh's quotient*. [E.13](#) When  $P_1$  is rank-one symmetric as in (2183),

$$\mathcal{R}(\text{vec } P_1 X P_1) = \mathcal{R}(\text{vec } P_1) \text{ in } \mathbb{R}^{m^2} \quad (2185)$$

and

$$P_1 X P_1 - X \perp P_1 \text{ in } \mathbb{R}^{m^2} \quad (2186)$$

The test for positive semidefiniteness, then, is a test for nonnegativity of the coefficient of orthogonal projection of  $X$  on the range of each and every vectorized extreme direction  $yy^T$  (§2.8.1) from the positive semidefinite cone in the ambient space of symmetric matrices.

#### E.6.4.3 $PXP \succeq 0$

In some circumstances, it may be desirable to limit the domain of test  $y^T X y \geq 0$  for positive semidefiniteness; e.g.,  $\{\|y\|=1\}$ . Another example of limiting domain-of-test is central to Euclidean distance geometry: For  $\mathcal{R}(V)=\mathcal{N}(\mathbf{1}^T)$ , the test  $-VDV \succeq 0$  determines whether  $D \in \mathbb{S}_h^N$  is a Euclidean distance matrix. The same test may be stated: For  $D \in \mathbb{S}_h^N$  (and optionally  $\|y\|=1$ )

$$D \in \mathbb{EDM}^N \Leftrightarrow -y^T D y = \langle yy^T, -D \rangle \geq 0 \quad \forall y \in \mathcal{R}(V) \quad (2187)$$

The test  $-VDV \succeq 0$  is therefore equivalent to a test for nonnegativity of the coefficient of orthogonal projection of  $-D$  on the range of each and every vectorized extreme direction  $yy^T$  from the positive semidefinite cone  $\mathbb{S}_+^N$  such that  $\mathcal{R}(yy^T) = \mathcal{R}(y) \subseteq \mathcal{R}(V)$ . (Validity of this result is independent of whether  $V$  is itself a projection matrix.)

---

[E.13](#)When  $y$  becomes the  $j^{\text{th}}$  eigenvector  $s_j$  of diagonalizable  $X$ , for example,  $\langle P_1, X \rangle$  becomes the  $j^{\text{th}}$  eigenvalue: [221, §III]

$$\langle P_1, X \rangle|_{y=s_j} = \frac{s_j^T \left( \sum_{i=1}^m \lambda_i s_i w_i^T \right) s_j}{s_j^T s_j} = \lambda_j$$

Similarly for  $y = w_j$ , the  $j^{\text{th}}$  left-eigenvector,

$$\langle P_1, X \rangle|_{y=w_j} = \frac{w_j^T \left( \sum_{i=1}^m \lambda_i s_i w_i^T \right) w_j}{w_j^T w_j} = \lambda_j$$

A quandary may arise regarding the potential annihilation of the antisymmetric part of  $X$  when  $s_j^T X s_j$  is formed. Were annihilation to occur, it would imply the eigenvalue thus found came instead from the symmetric part of  $X$ . The quandary is resolved recognizing that diagonalization of real  $X$  admits complex eigenvectors; hence, annihilation could only come about by forming  $\text{re}(s_j^H X s_j) = s_j^H (X + X^T) s_j / 2$  [228, §7.1] where  $(X + X^T)/2$  is the symmetric part of  $X$ , and  $s_j^H$  denotes conjugate transpose.

## E.7 Projection on matrix subspaces

### E.7.1 $PXP$ interpretation for higher-rank $P$

For a projection matrix  $P$  of rank greater than 1,  $PXP$  is generally not commensurate with  $\frac{\langle P, X \rangle}{\langle P, P \rangle}P$  as is the case for projector dyads (2184). Yet for a symmetric idempotent matrix  $P$  of any rank we are tempted to say, erroneously, “ $PXP$  is the orthogonal projection of  $X \in \mathbb{S}^m$  on  $\mathcal{R}(\text{vec } P)$ ”. The fallacy is:  $\text{vec } PXP$  does not necessarily belong to the range of vectorized  $P$ ; the most basic requirement for projection on  $\mathcal{R}(\text{vec } P)$ .

**E.7.1.0.1 Theorem.** *Kronecker projector.* [363, §2.7]

Given any projection matrices  $P_1$  and  $P_2$  (subspace projectors), then

$$P_1 \otimes P_2, \quad P_1 \otimes I \quad (2188)$$

are projection matrices. The product preserves symmetry when present.  $\diamond$

But for  $P$  of any rank, by this theorem, we may always say:  $PXP$  is the orthogonal projection of  $X \in \mathbb{S}^m$  on  $\mathcal{R}(P \otimes P)$  because  $\text{vec } PXP = (P \otimes P) \text{vec } X$  (§A.1.1 no.33). Only when projection matrix  $P$  has rank 1 may we say

$$\begin{aligned} \text{vec } PXP &= (P \otimes P) \text{vec } X = \langle P, X \rangle \text{vec } P \\ \mathcal{R}(\text{vec } PXP) &= \mathcal{R}(P \otimes P) = \mathcal{R}(\text{vec } P) \end{aligned} \quad (2189)$$

### E.7.2 Orthogonal projection on matrix subspaces

With  $A_1 \in \mathbb{R}^{m \times n}$ ,  $B_1 \in \mathbb{R}^{n \times k}$ ,  $Z_1 \in \mathbb{R}^{m \times k}$ ,  $A_2 \in \mathbb{R}^{p \times n}$ ,  $B_2 \in \mathbb{R}^{n \times k}$ ,  $Z_2 \in \mathbb{R}^{p \times k}$  as defined for nonorthogonal projector (2069), and defining

$$P_1 \triangleq A_1 A_1^\dagger \in \mathbb{S}^m, \quad P_2 \triangleq A_2 A_2^\dagger \in \mathbb{S}^p \quad (2190)$$

then, given dimensionally compatible  $X$

$$\|X - P_1 X P_2\|_F = \inf_{B_1, B_2 \in \mathbb{R}^{n \times k}} \|X - A_1(A_1^\dagger + B_1 Z_1^\top) X (A_2^{\dagger T} + Z_2 B_2^\top) A_2^T\|_F \quad (2191)$$

As for all subspace projectors, range of the projector is the subspace on which projection is made:  $\{P_1 Y P_2 \mid Y \in \mathbb{R}^{m \times p}\}$ . For projectors  $P_1$  and  $P_2$  of any rank, altogether, this means projection  $P_1 X P_2$  is unique minimum-distance, orthogonal

$$P_1 X P_2 - X \perp \{P_1 Y P_2 \mid Y \in \mathbb{R}^{m \times p}\} \text{ in } \mathbb{R}^{mp} \quad (2192)$$

and  $P_1$  and  $P_2$  must each be symmetric (*confer*(2172)) to attain the infimum.

**E.7.2.0.1 Proof.** *Minimum Frobenius norm (2191).*

Defining  $P \triangleq A_1(A_1^\dagger + B_1 Z_1^\top)$ ,

$$\begin{aligned} &\inf_{B_1, B_2} \|X - A_1(A_1^\dagger + B_1 Z_1^\top) X (A_2^{\dagger T} + Z_2 B_2^\top) A_2^T\|_F^2 \\ &= \inf_{B_1, B_2} \|X - P X (A_2^{\dagger T} + Z_2 B_2^\top) A_2^T\|_F^2 \\ &= \inf_{B_1, B_2} \text{tr}((X^\top - A_2(A_2^\dagger + B_2 Z_2^\top) X^\top P^\top)(X - P X (A_2^{\dagger T} + Z_2 B_2^\top) A_2^T)) \\ &= \inf_{B_1, B_2} \text{tr}(X^\top X - X^\top P X (A_2^{\dagger T} + Z_2 B_2^\top) A_2^T - A_2(A_2^\dagger + B_2 Z_2^\top) X^\top P^\top X \\ &\quad + A_2(A_2^\dagger + B_2 Z_2^\top) X^\top P^\top P X (A_2^{\dagger T} + Z_2 B_2^\top) A_2^T) \end{aligned} \quad (2193)$$

Necessary conditions for a global minimum are  $\nabla_{B_1} = \mathbf{0}$  and  $\nabla_{B_2} = \mathbf{0}$ . Terms not containing  $B_2$  in (2193) will vanish from gradient  $\nabla_{B_2}$ ; (§D.2.3)

$$\begin{aligned} \nabla_{B_2} \text{tr} & \left( -X^T P X Z_2 B_2 A_2^T - A_2 B_2 Z_2^T X^T P^T X + A_2 A_2^\dagger X^T P^T P X Z_2 B_2^T A_2^T \right. \\ & \quad \left. + A_2 B_2 Z_2^T X^T P^T P X A_2^{\dagger T} A_2^T + A_2 B_2 Z_2^T X^T P^T P X Z_2 B_2^T A_2^T \right) \\ &= -2A_2^T X^T P X Z_2 + 2A_2^T A_2^\dagger X^T P^T P X Z_2 + 2A_2^T A_2 B_2 Z_2^T X^T P^T P X Z_2 \\ &= A_2^T \left( -X^T + A_2 A_2^\dagger X^T P^T + A_2 B_2 Z_2^T X^T P^T \right) P X Z_2 \\ &= \mathbf{0} \qquad \Leftrightarrow \\ \mathcal{R}(B_1) &\subseteq \mathcal{N}(A_1) \quad \text{and} \quad \mathcal{R}(B_2) \subseteq \mathcal{N}(A_2) \end{aligned} \tag{2194}$$

(or  $Z_2 = \mathbf{0}$ ) because  $A^T = A^T A A^\dagger$ . Symmetry requirement (2190) is implicit. Were instead  $P^T \triangleq (A_2^{\dagger T} + Z_2 B_2^T) A_2^T$  and the gradient with respect to  $B_1$  observed, then similar results are obtained. The projector is unique. Perpendicularity (2192) establishes uniqueness [122, §4.9] of projection  $P_1 X P_2$  on a matrix subspace. The minimum-distance projector is the orthogonal projector, and *vice versa*. ♦

#### E.7.2.0.2 Example. $PXP$ redux & $\mathcal{N}(\mathbf{V})$ .

Suppose we define a subspace of  $m \times n$  matrices, each elemental matrix having columns constituting a list whose geometric center (§5.5.1.0.1) is the origin in  $\mathbb{R}^m$ :

$$\begin{aligned} \mathbb{R}_c^{m \times n} &\triangleq \{Y \in \mathbb{R}^{m \times n} \mid Y\mathbf{1} = \mathbf{0}\} \\ &= \{Y \in \mathbb{R}^{m \times n} \mid \mathcal{N}(Y) \supseteq \mathbf{1}\} = \{Y \in \mathbb{R}^{m \times n} \mid \mathcal{R}(Y^T) \subseteq \mathcal{N}(\mathbf{1}^T)\} \\ &= \{XV \mid X \in \mathbb{R}^{m \times n}\} \subset \mathbb{R}^{m \times n} \end{aligned} \tag{2195}$$

the *nonsymmetric geometric center subspace*. Further suppose  $V \in \mathbb{S}^n$  is a projection matrix having  $\mathcal{N}(V) = \mathcal{R}(\mathbf{1})$  and  $\mathcal{R}(V) = \mathcal{N}(\mathbf{1}^T)$ . Then linear mapping  $T(X) = XV$  is the orthogonal projection of any  $X \in \mathbb{R}^{m \times n}$  on  $\mathbb{R}_c^{m \times n}$  in the Euclidean (vectorization) sense because  $V$  is symmetric,  $\mathcal{N}(XV) \supseteq \mathbf{1}$ , and  $\mathcal{R}(VX^T) \subseteq \mathcal{N}(\mathbf{1}^T)$ .

Now suppose we define a subspace of symmetric  $n \times n$  matrices each of whose columns constitute a list having the origin in  $\mathbb{R}^n$  as geometric center,

$$\begin{aligned} \mathbb{S}_c^n &\triangleq \{Y \in \mathbb{S}^n \mid Y\mathbf{1} = \mathbf{0}\} \\ &= \{Y \in \mathbb{S}^n \mid \mathcal{N}(Y) \supseteq \mathbf{1}\} = \{Y \in \mathbb{S}^n \mid \mathcal{R}(Y) \subseteq \mathcal{N}(\mathbf{1}^T)\} \end{aligned} \tag{2196}$$

the *geometric center subspace*. Further suppose  $V \in \mathbb{S}^n$  is a projection matrix, the same as before. Then  $VXV$  is the orthogonal projection of any  $X \in \mathbb{S}^n$  on  $\mathbb{S}_c^n$  in the Euclidean sense (2192) because  $V$  is symmetric,  $VXV\mathbf{1} = \mathbf{0}$ , and  $\mathcal{R}(VXV) \subseteq \mathcal{N}(\mathbf{1}^T)$ . Two-sided projection is necessary only to remain in the ambient symmetric matrix subspace. Then

$$\mathbb{S}_c^n = \{VXV \mid X \in \mathbb{S}^n\} \subset \mathbb{S}^n \tag{2197}$$

has  $\dim \mathbb{S}_c^n = n(n-1)/2$  in isomorphic  $\mathbb{R}^{n(n+1)/2}$ . We find its orthogonal complement as the aggregate of all negative directions of orthogonal projection on  $\mathbb{S}_c^n$ : the *translation-invariant subspace* (§5.5.1.1)

$$\begin{aligned} \mathbb{S}_c^{n \perp} &\triangleq \{X - VXV \mid X \in \mathbb{S}^n\} \subset \mathbb{S}^n \\ &= \{u\mathbf{1}^T + \mathbf{1}u^T \mid u \in \mathbb{R}^n\} \end{aligned} \tag{2198}$$

characterized by doublet  $u\mathbf{1}^T + \mathbf{1}u^T$  (§B.2). [E.14](#) Defining geometric center mapping

$$\mathbf{V}(X) = -\frac{1}{2}VXV \quad (1141)$$

consistently with (1141), then  $\mathcal{N}(\mathbf{V}) = \mathcal{R}(I - \mathbf{V})$  on domain  $\mathbb{S}^n$  analogously to vector projectors (§E.2); *id est*,

$$\mathcal{N}(\mathbf{V}) = \mathbb{S}_c^{n\perp} \quad (2199)$$

a subspace of  $\mathbb{S}^n$  whose dimension is  $\dim \mathbb{S}_c^{n\perp} = n$  in isomorphic  $\mathbb{R}^{n(n+1)/2}$ . Intuitively, operator  $\mathbf{V}$  is an orthogonal projector; any argument duplicitously in its range is a fixed point. So, this symmetric operator's nullspace must be orthogonal to its range.

Now compare the subspace of symmetric matrices having all zeros in the first row and column

$$\begin{aligned} \mathbb{S}_{\mathbf{0}}^n &\triangleq \{Y \in \mathbb{S}^n \mid Ye_1 = \mathbf{0}\} \\ &= \left\{ \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} X \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} \mid X \in \mathbb{S}^n \right\} \\ &= \left\{ [\mathbf{0} \ \sqrt{2}V_{\mathcal{N}}]^T Z [\mathbf{0} \ \sqrt{2}V_{\mathcal{N}}] \mid Z \in \mathbb{S}^n \right\} \end{aligned} \quad (2200)$$

where  $P = \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix}$  is an orthogonal projector and  $[\mathbf{0} \ \sqrt{2}V_{\mathcal{N}}]$  is a nonorthogonal projector (§B.4.2 *no.7*). Then, similarly,  $PXP$  is the orthogonal projection of any  $X \in \mathbb{S}^n$  on  $\mathbb{S}_{\mathbf{0}}^n$  in the Euclidean sense (2192). Like  $\mathbb{S}_c^n$  (§6.8.1.1.1),  $\mathbb{S}_{\mathbf{0}}^n$  is invariant to projection on a positive semidefinite cone. The orthogonal complement of  $\mathbb{S}_{\mathbf{0}}^n$  is

$$\begin{aligned} \mathbb{S}_{\mathbf{0}}^{n\perp} &\triangleq \left\{ \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} X \begin{bmatrix} \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & I \end{bmatrix} - X \mid X \in \mathbb{S}^n \right\} \subset \mathbb{S}^n \\ &= \{ue_1^T + e_1u^T \mid u \in \mathbb{R}^n\} \end{aligned} \quad (2201)$$

Obviously,  $\mathbb{S}_{\mathbf{0}}^n \oplus \mathbb{S}_{\mathbf{0}}^{n\perp} = \mathbb{S}^n$ . □

## E.8 Range Rowspace interpretation

For idempotent matrices  $P_1$  and  $P_2$  of any rank,  $P_1XP_2^T$  is a projection of  $\mathcal{R}(X)$  on  $\mathcal{R}(P_1)$  and a projection of  $\mathcal{R}(X^T)$  on  $\mathcal{R}(P_2)$ : For any given  $X = U\Sigma Q^T = \sum_{i=1}^{\eta} \sigma_i u_i q_i^T \in \mathbb{R}^{m \times p}$ , as in compact SVD (1716),

$$P_1XP_2^T = \sum_{i=1}^{\eta} \sigma_i P_1 u_i q_i^T P_2^T = \sum_{i=1}^{\eta} \sigma_i P_1 u_i (P_2 q_i)^T \quad (2202)$$

where  $\eta \triangleq \min\{m, p\}$ . Recall:  $u_i \in \mathcal{R}(X)$  and  $q_i \in \mathcal{R}(X^T)$  when the corresponding singular value  $\sigma_i$  is nonzero. (§A.6.1) So  $P_1$  projects  $u_i$  on  $\mathcal{R}(P_1)$  while  $P_2$  projects

---

**E.14 Proof.**

$$\begin{aligned} \{X - VXV \mid X \in \mathbb{S}^n\} &= \{X - (I - \frac{1}{n}\mathbf{1}\mathbf{1}^T)X(I - \mathbf{1}\mathbf{1}^T\frac{1}{n}) \mid X \in \mathbb{S}^n\} \\ &= \{\frac{1}{n}\mathbf{1}\mathbf{1}^T X + X\mathbf{1}\mathbf{1}^T\frac{1}{n} - \frac{1}{n}\mathbf{1}\mathbf{1}^T X \mathbf{1}\mathbf{1}^T\frac{1}{n} \mid X \in \mathbb{S}^n\} \end{aligned}$$

Because  $\{X\mathbf{1} \mid X \in \mathbb{S}^n\} = \mathbb{R}^n$ ,

$$\begin{aligned} \{X - VXV \mid X \in \mathbb{S}^n\} &= \{\mathbf{1}\zeta^T + \zeta\mathbf{1}^T - \mathbf{1}\mathbf{1}^T(\mathbf{1}^T\zeta\frac{1}{n}) \mid \zeta \in \mathbb{R}^n\} \\ &= \{\mathbf{1}\zeta^T(I - \mathbf{1}\mathbf{1}^T\frac{1}{2n}) + (I - \frac{1}{2n}\mathbf{1}\mathbf{1}^T)\zeta\mathbf{1}^T \mid \zeta \in \mathbb{R}^n\} \end{aligned}$$

where  $I - \frac{1}{2n}\mathbf{1}\mathbf{1}^T$  is invertible. ♦

$q_i$  on  $\mathcal{R}(P_2)$ ; *id est*, the range and rowspace of any  $X$  are respectively projected on the ranges of  $P_1$  and  $P_2$ .<sup>E.15</sup>

## E.9 Projection on convex set

Thus far we have discussed only projection on subspaces. Now we generalize, considering projection on arbitrary convex sets in Euclidean space; convex because point of projection is then unique, minimum-distance, and a convex optimization problem:

For projection  $P_C x$  of point  $x$  on any closed set  $C \subseteq \mathbb{R}^n$  it is obvious:

$$\mathcal{C} \equiv \{P_C x \mid x \in \mathbb{R}^n\} = \{x \in \mathbb{R}^n \mid P_C x = x\} \quad (2203)$$

where  $P_C$  is a projection operator that is convex when  $C$  is convex. [65, p.88]

If  $C \subseteq \mathbb{R}^n$  is a closed convex set, then for each and every  $x \in \mathbb{R}^n$  there exists a unique point  $P_C x$  belonging to  $C$  that is closest to  $x$  in the Euclidean sense. Like (2091), unique projection  $Px$  (or  $P_C x$ ) of a point  $x$  on convex set  $C$  is that point in  $C$  closest to  $x$ ; [280, §3.12]

$$\|x - Px\|_2 = \inf_{y \in C} \|x - y\|_2 = \text{dist}(x, C) \quad (2204)$$

There exists a converse (in finite-dimensional Euclidean space):

**E.9.0.0.1 Theorem.** (Bunt-Motzkin) *Convex set if projections unique.* [434, §7.5] [222] If  $C \subseteq \mathbb{R}^n$  is a nonempty closed set and if for each and every  $x$  in  $\mathbb{R}^n$  there is a unique Euclidean projection  $Px$  of  $x$  on  $C$  belonging to  $C$ , then  $C$  is convex.  $\diamond$

Borwein & Lewis propose, for closed convex set  $C$  [58, §3.3 exer.12d]

$$\nabla \|x - Px\|_2^2 = 2(x - Px) \quad (2205)$$

for any point  $x$  whereas, for  $x \notin C$

$$\nabla \|x - Px\|_2 = (x - Px) \|x - Px\|_2^{-1} \quad (2206)$$

**E.9.0.0.2 Exercise.** *Norm gradient.*

Prove (2205) and (2206). (Not proved in [58].)  $\blacktriangledown$

A well-known equivalent characterization of projection on a convex set is a generalization of the perpendicularity condition (2090) for projection on a subspace:

### E.9.1 Dual interpretation of projection on convex set

**E.9.1.0.1 Definition.** *Normal vector.*

[343, p.15]

Vector  $z$  is *normal* to convex set  $C$  at point  $Px \in C$  if

$$\langle z, y - Px \rangle \leq 0 \quad \forall y \in C \quad (2207)$$

$\triangle$

A convex set has at least one nonzero normal at each of its boundary points. [343, p.100] (Figure 71) Hence, the *normal* or *dual* interpretation of projection:

---

<sup>E.15</sup>When  $P_1$  and  $P_2$  are symmetric and  $\mathcal{R}(P_1) = \mathcal{R}(u_j)$  and  $\mathcal{R}(P_2) = \mathcal{R}(q_j)$ , then the  $j^{\text{th}}$  dyad term from the singular value decomposition of  $X$  is isolated by the projection. Yet if  $\mathcal{R}(P_2) = \mathcal{R}(q_\ell)$ ,  $\ell \neq j \in \{1 \dots \eta\}$ , then  $P_1 X P_2 = \mathbf{0}$ .

**E.9.1.0.2 Theorem.** *Unique minimum-distance projection.* [225, §A.3.1] [280, §3.12] [122, §4.1] [86] (Figure 197b p.609) Given a closed convex set  $\mathcal{C} \subseteq \mathbb{R}^n$ , point  $Px$  is the unique projection of a given point  $x \in \mathbb{R}^n$  on  $\mathcal{C}$  ( $Px$  is that point in  $\mathcal{C}$  nearest  $x$ ) if and only if

$$Px \in \mathcal{C}, \quad \langle x - Px, y - Px \rangle \leq 0 \quad \forall y \in \mathcal{C} \quad (2208)$$

◊

As for subspace projection, convex operator  $P$  is idempotent in the sense: for each and every  $x \in \mathbb{R}^n$ ,  $P(Px) = Px$ . Yet operator  $P$  is nonlinear;

- Projector  $P$  is a linear operator if and only if convex set  $\mathcal{C}$  (on which projection is made) is a subspace. (§E.4)

**E.9.1.0.3 Theorem.** *Unique projection via normal cone.* [E.16](#) [122, §4.3] Given closed convex set  $\mathcal{C} \subseteq \mathbb{R}^n$ , point  $Px$  is the unique projection of a given point  $x \in \mathbb{R}^n$  on  $\mathcal{C}$  if and only if

$$Px \in \mathcal{C}, \quad Px - x \in (\mathcal{C} - Px)^* \quad (2209)$$

In other words,  $Px$  is that point in  $\mathcal{C}$  nearest  $x$  if and only if  $Px - x$  belongs to that cone dual to translate  $\mathcal{C} - Px$ . ◊

### E.9.1.1 Dual interpretation as optimization

Deutsch [124, thm.2.3] [125, §2] and Luenberger [280, p.134] carry forward Nirenberg's dual interpretation of projection [311] as solution to a maximization problem: Minimum distance from a point  $x \in \mathbb{R}^n$  to a convex set  $\mathcal{C} \subset \mathbb{R}^n$  can be found by maximizing distance from  $x$  to hyperplane  $\partial\mathcal{H}$  over the set of all hyperplanes separating  $x$  from  $\mathcal{C}$ . Existence of a separating hyperplane (§2.4.2.7) presumes that point  $x$  lies on the boundary or exterior to set  $\mathcal{C}$ .

The optimal separating hyperplane is characterized by the fact that it also supports  $\mathcal{C}$ . Any hyperplane supporting  $\mathcal{C}$  (Figure 32a) has form

$$\underline{\partial\mathcal{H}} = \{y \in \mathbb{R}^n \mid a^T y = \sigma_{\mathcal{C}}(a)\} \quad (131)$$

where support function

$$\sigma_{\mathcal{C}}(a) = \sup_{z \in \mathcal{C}} a^T z \quad (560)$$

is convex w.r.t  $a$ . When point  $x$  is finite and set  $\mathcal{C}$  contains finite points, under this dual interpretation, if the supporting hyperplane is a separating hyperplane then the support function is finite. From Example E.5.0.0.5, projection  $P_{\underline{\partial\mathcal{H}}} x$  of  $x$  on any given supporting hyperplane  $\underline{\partial\mathcal{H}}$  is

$$P_{\underline{\partial\mathcal{H}}} x = x - a(a^T a)^{-1}(a^T x - \sigma_{\mathcal{C}}(a)) \quad (2210)$$

With reference to Figure 191, identifying

$$\mathcal{H}_+ = \{y \in \mathbb{R}^n \mid a^T y \geq \sigma_{\mathcal{C}}(a)\} \quad (109)$$

then

$$\begin{aligned} \|x - P_{\mathcal{C}} x\| &= \sup_{\underline{\partial\mathcal{H}} \mid x \in \mathcal{H}_+} \|x - P_{\underline{\partial\mathcal{H}}} x\| = \sup_{a \mid x \in \mathcal{H}_+} \|a(a^T a)^{-1}(a^T x - \sigma_{\mathcal{C}}(a))\| \\ &= \sup_{a \mid x \in \mathcal{H}_+} \frac{1}{\|a\|} |a^T x - \sigma_{\mathcal{C}}(a)| \end{aligned} \quad (2211)$$

---

[E.16](#)  $-(\mathcal{C} - Px)^*$  is the normal cone to set  $\mathcal{C}$  at point  $Px$ . (§E.10.3.2)

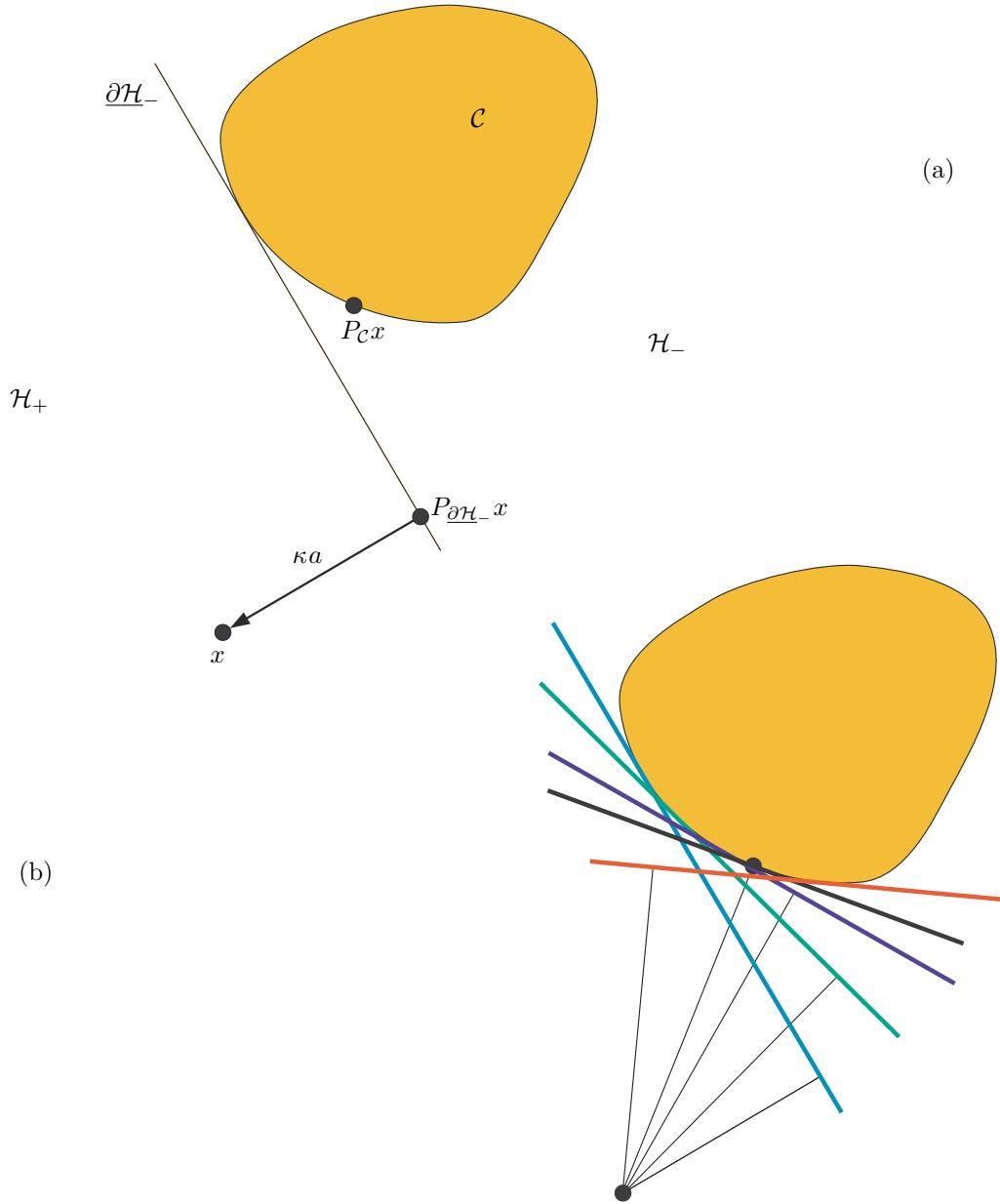


Figure 191: Dual interpretation of projection of point  $x$  on convex set  $\mathcal{C}$  in  $\mathbb{R}^2$ .  
**(a)**  $\kappa = (a^T a)^{-1} (a^T x - \sigma_{\mathcal{C}}(a))$ .    **(b)** Minimum distance from  $x$  to  $\mathcal{C}$  is found by maximizing distance to all hyperplanes supporting  $\mathcal{C}$  and separating it from  $x$ . Distance of maximization is unique over any convex set.

which can be expressed, for arbitrary positive constant  $\tau$

$$\begin{aligned} \|x - P_{\mathcal{C}}x\| &= \frac{1}{\tau} \underset{a}{\text{maximize}} \quad a^T x - \sigma_{\mathcal{C}}(a) \\ &\text{subject to} \quad \|a\| \leq \tau \end{aligned} \quad (2212)$$

The unique minimum-distance projection on convex set  $\mathcal{C}$  is therefore

$$P_{\mathcal{C}}x = x - a^* (a^{*T}x - \sigma_{\mathcal{C}}(a^*)) \frac{1}{\tau^2} \quad (2213)$$

where optimally  $\|a^*\| = \tau$ .

#### E.9.1.1.1 Exercise. Dual projection technique on polyhedron.

Test that projection paradigm from Figure 191 on any convex polyhedral set. ▼

#### E.9.1.1.2 Exercise. Projection on boundary from inside.

Now suppose that point  $x$  lies interior to convex set  $\mathcal{C}$ . What is the consequence of

$$\begin{aligned} \frac{1}{\tau} \underset{a}{\text{minimize}} \quad &\sigma_{\mathcal{C}}(a) - a^T x \\ &\text{subject to} \quad \|a\| = \tau \end{aligned} \quad (2214)$$

Is this program convex?<sup>E.17</sup> Why can we not say  $\|a\| \leq \tau$  here? State conditions under which a boundary solution is unique. ▼

#### E.9.1.2 Dual interpretation of projection on cone

In the circumstance that set  $\mathcal{C}$  is a closed convex cone  $\mathcal{K}$  and there exists a hyperplane separating given point  $x$  from  $\mathcal{K}$ , then optimal  $\sigma_{\mathcal{K}}(a^*)$  takes value 0 [225, §C.2.3.1]. So problem (2212) for projection of  $x$  on  $\mathcal{C} = \mathcal{K}$  becomes convex:

$$\begin{aligned} \|x - P_{\mathcal{K}}x\| &= \frac{1}{\tau} \underset{a}{\text{maximize}} \quad a^T x \\ &\text{subject to} \quad \|a\| \leq \tau \\ &\quad a \in \mathcal{K}^\circ \end{aligned} \quad (2215)$$

Here, the norm inequality can be handled by Schur complement (§3.5.3). Normals  $a$  to all hyperplanes supporting  $\mathcal{K}$  belong to the polar cone  $\mathcal{K}^\circ = -\mathcal{K}^*$  by definition: (322)

$$a \in \mathcal{K}^\circ \Leftrightarrow \langle a, x \rangle \leq 0 \quad \text{for all } x \in \mathcal{K} \quad (2216)$$

Projection on cone  $\mathcal{K}$  is

$$P_{\mathcal{K}}x = (I - \frac{1}{\tau^2} a^* a^{*T})x \quad (2217)$$

whereas projection on the polar cone  $-\mathcal{K}^*$  is (§E.9.2.2.1)

$$P_{\mathcal{K}^\circ}x = x - P_{\mathcal{K}}x = \frac{1}{\tau^2} a^* a^{*T} x \quad (2218)$$

Negating vector  $a$ , this convex maximization problem (2215) becomes a minimization (the same problem) and the polar cone becomes the dual cone:

$$\begin{aligned} \|x - P_{\mathcal{K}}x\| &= -\frac{1}{\tau} \underset{a}{\text{minimize}} \quad a^T x \\ &\text{subject to} \quad \|a\| \leq \tau \\ &\quad a \in \mathcal{K}^* \end{aligned} \quad (2219)$$

---

<sup>E.17</sup> Hint: §4.7.0.0.1.

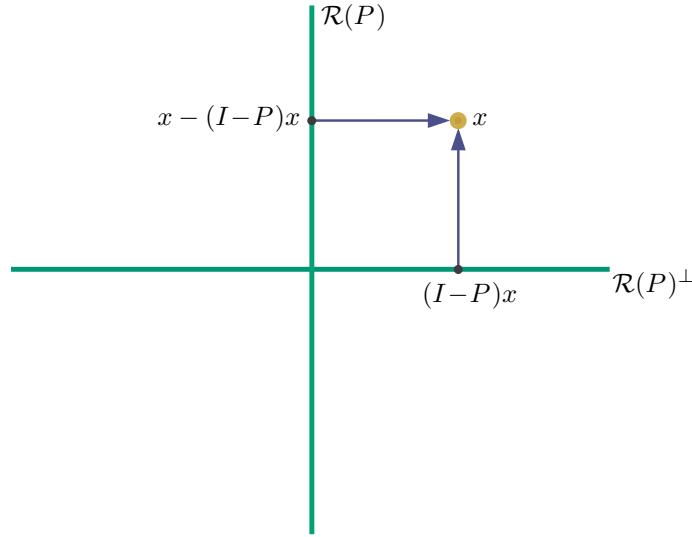


Figure 192: (confer Figure 95, Figure 193) Given orthogonal projection  $(I-P)x$  of  $x$  on orthogonal complement  $\mathcal{R}(P)^\perp$ , projection on  $\mathcal{R}(P)$  is immediate:  $x - (I-P)x$ .

### E.9.2 Projection on cone

When convex set  $\mathcal{C}$  is a cone, there is a finer statement of optimality conditions:

**E.9.2.0.1 Theorem.** *Unique projection on cone.* [225, §A.3.2]

Let  $\mathcal{K} \subseteq \mathbb{R}^n$  be a closed convex cone, and  $\mathcal{K}^*$  its dual (§2.13.1). Then  $Px$  is the unique minimum-distance projection of  $x \in \mathbb{R}^n$  on  $\mathcal{K}$  if and only if

$$Px \in \mathcal{K}, \quad \langle Px - x, Px \rangle = 0, \quad Px - x \in \mathcal{K}^* \quad (2220)$$

◊

In words,  $Px$  is the unique minimum-distance projection of  $x$  on  $\mathcal{K}$  if and only if

- 1) projection  $Px$  lies in  $\mathcal{K}$
- 2) direction  $Px - x$  is orthogonal to the projection  $Px$
- 3) direction  $Px - x$  lies in the dual cone  $\mathcal{K}^*$ .

As the theorem is stated, it admits projection on  $\mathcal{K}$  not full-dimensional; *id est*, on closed convex cones in a proper subspace of  $\mathbb{R}^n$ .

Projection on  $\mathcal{K}$  of any point  $x \in -\mathcal{K}^*$ , belonging to the negative dual cone, is the origin. By (2220): the set of all points reaching the origin, when projecting on  $\mathcal{K}$ , constitutes the negative dual cone; **a.k.a**, the *polar cone*

$$\mathcal{K}^\circ = -\mathcal{K}^* = \{x \in \mathbb{R}^n \mid Px = \mathbf{0}\} \quad (2221)$$

#### E.9.2.1 Relationship to subspace projection

Conditions 1 and 2 of Theorem E.9.2.0.1 are common with orthogonal projection on a subspace  $\mathcal{R}(P)$ :

- 1) Condition 1 corresponds to the most basic requirement; namely, the projection  $Px \in \mathcal{R}(P)$  belongs to the subspace (*confer*(2203))

$$\mathcal{K} = \{Px \mid x \in \mathbb{R}^n\} \triangleq \mathcal{R}(P) \quad (2222)$$

- 2) Recall the perpendicularity requirement for projection on a subspace:

$$Px - x \perp \mathcal{R}(P) \quad \text{or} \quad Px - x \in \mathcal{R}(P)^\perp \quad (2090)$$

which corresponds to condition 2.

- 3) Yet condition 3 is a generalization of subspace projection; *id est*, for unique minimum-distance projection on a closed convex cone  $\mathcal{K}$ , polar cone  $-\mathcal{K}^*$  (Figure 193) plays the role that  $\mathcal{R}(P)^\perp$  plays for subspace projection (Figure 192):

$$P_{\mathcal{R}}x = x - P_{\mathcal{R}^\perp}x \quad (2223)$$

Indeed,  $-\mathcal{K}^*$  is the algebraic complement in the orthogonal vector sum (p.624) [300] [225, §A.3.2.5]

$$\mathcal{K} \boxplus -\mathcal{K}^* = \mathbb{R}^n \Leftrightarrow \text{cone } \mathcal{K} \text{ is closed and convex} \quad (2224)$$

Given unique minimum-distance projection  $Px$  on  $\mathcal{K}$  satisfying Theorem E.9.2.0.1, then by projection on the algebraic complement via  $I - P$  in §E.2 we have

$$-\mathcal{K}^* = \{x - Px \mid x \in \mathbb{R}^n\} = \{x \in \mathbb{R}^n \mid Px = \mathbf{0}\} = \mathcal{N}(P) \quad (2225)$$

consequent to Moreau (2228). Converse (2225)(2222)  $\Rightarrow$  (2224) holds as well. Recalling that any subspace is a closed convex cone<sup>E.18</sup>

$$\mathcal{K} = \mathcal{R}(P) \Leftrightarrow -\mathcal{K}^* = \mathcal{R}(P)^\perp \quad (2226)$$

meaning, when a cone is a subspace  $\mathcal{R}(P)$ , then the dual cone becomes its orthogonal complement  $\mathcal{R}(P)^\perp$ . In this circumstance, condition 3 becomes coincident with condition 2.

Properties, of projection on cones in what follows, further generalize to subspaces by: (4)

$$\mathcal{K} = \mathcal{R}(P) \Leftrightarrow -\mathcal{K} = \mathcal{R}(P) \quad (2227)$$

### E.9.2.2 Salient properties: Projection $Px$ on closed convex cone $\mathcal{K}$

[225, §A.3.2] [122, §5.6] For  $x, x_1, x_2 \in \mathbb{R}^n$

1.  $P_{\mathcal{K}}(\alpha x) = \alpha P_{\mathcal{K}}x \quad \forall \alpha \geq 0$  (nonnegative homogeneity)
2.  $\|P_{\mathcal{K}}x\| \leq \|x\|$
3.  $P_{\mathcal{K}}x = \mathbf{0} \Leftrightarrow x \in -\mathcal{K}^*$
4.  $P_{\mathcal{K}}(-x) = -P_{-\mathcal{K}}x$
5. (Jean-Jacques Moreau, 1962) [300]

$$\begin{aligned} x &= x_1 + x_2, \quad x_1 \in \mathcal{K}, \quad x_2 \in -\mathcal{K}^*, \quad x_1 \perp x_2 \\ &\Leftrightarrow \\ x_1 &= P_{\mathcal{K}}x, \quad x_2 = P_{-\mathcal{K}^*}x \end{aligned} \quad (2228)$$

6.  $\mathcal{K} = \{x - P_{-\mathcal{K}^*}x \mid x \in \mathbb{R}^n\} = \{x \in \mathbb{R}^n \mid P_{-\mathcal{K}^*}x = \mathbf{0}\}$
7.  $-\mathcal{K}^* = \{x - P_{\mathcal{K}}x \mid x \in \mathbb{R}^n\} = \{x \in \mathbb{R}^n \mid P_{\mathcal{K}}x = \mathbf{0}\}$  (2225)

---

<sup>E.18</sup> but a proper subspace is not a proper cone (§2.7.2.2.1).

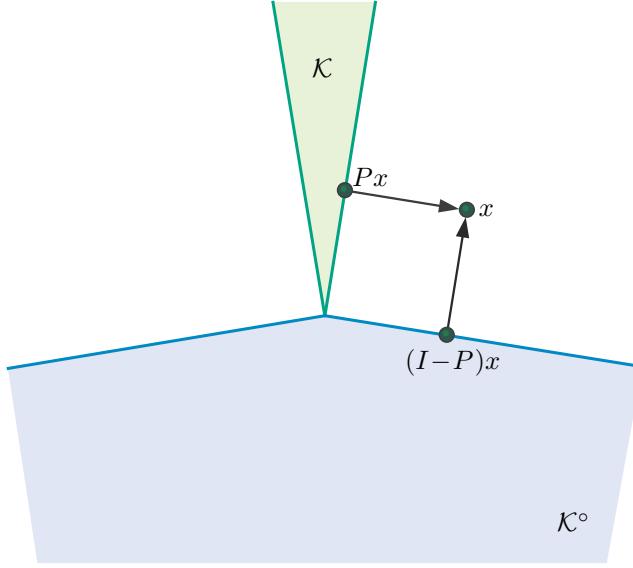


Figure 193: (confer Figure 192) Given minimum-distance projection  $(I-P)x$  of  $x$  on negative dual cone  $\mathcal{K}^\circ$ , projection on  $\mathcal{K}$  is immediate:  $x - (I-P)x = Px$ .

**E.9.2.2.1 Corollary.**  $I - P$  for cones. (confer §E.2)

Denote by  $\mathcal{K} \subseteq \mathbb{R}^n$  a closed convex cone, and call  $\mathcal{K}^*$  its dual. Then  $x - P_{\mathcal{K}^*}x$  is the unique minimum-distance projection of  $x \in \mathbb{R}^n$  on  $\mathcal{K}$  if and only if  $P_{\mathcal{K}^*}x$  is the unique minimum-distance projection of  $x$  on  $-\mathcal{K}^*$  the polar cone. ◊

**Proof.** Assume  $x_1 = P_{\mathcal{K}}x$ . Then by Theorem E.9.2.0.1 we have

$$x_1 \in \mathcal{K}, \quad x_1 - x \perp x_1, \quad x_1 - x \in \mathcal{K}^* \quad (2229)$$

Now assume  $x - x_1 = P_{\mathcal{K}^*}x$ . Then we have

$$x - x_1 \in -\mathcal{K}^*, \quad -x_1 \perp x - x_1, \quad -x_1 \in -\mathcal{K} \quad (2230)$$

But these two assumptions are apparently identical. We must therefore have

$$x - P_{\mathcal{K}^*}x = x_1 = P_{\mathcal{K}}x \quad (2231)$$

♦

**E.9.2.2.2 Corollary.** Unique projection via dual or normal cone. [122, §4.7]

(§E.10.3.2, confer Theorem E.9.1.0.3) Given point  $x \in \mathbb{R}^n$  and closed convex cone  $\mathcal{K} \subseteq \mathbb{R}^n$ , the following are equivalent statements:

1. point  $Px$  is the unique minimum-distance projection of  $x$  on  $\mathcal{K}$
2.  $Px \in \mathcal{K}$ ,  $x - Px \in -(\mathcal{K} - Px)^* = -\mathcal{K}^* \cap (Px)^\perp$
3.  $Px \in \mathcal{K}$ ,  $\langle x - Px, Px \rangle = 0$ ,  $\langle x - Px, y \rangle \leq 0 \quad \forall y \in \mathcal{K}$  ◊

**E.9.2.2.3 Example.** *Unique projection on nonnegative orthant.* (confer (1466))

To project matrix  $H \in \mathbb{R}^{m \times n}$  on the selfdual orthant ([§2.13.6.1](#)) of nonnegative matrices  $\mathbb{R}_+^{m \times n}$  in isomorphic  $\mathbb{R}^{mn}$ , from Theorem [E.9.2.0.1](#), necessary and sufficient conditions are:

$$\begin{aligned} H^* &\geq \mathbf{0} \\ \text{tr}((H^* - H)^T H^*) &= 0 \\ H^* - H &\geq \mathbf{0} \end{aligned} \tag{2232}$$

where the inequalities denote entrywise comparison. The optimal solution  $H^*$  is simply  $H$  having all its negative entries zeroed;

$$H_{ij}^* = \max\{H_{ij}, 0\}, \quad i, j \in \{1 \dots m\} \times \{1 \dots n\} \tag{2233}$$

Now suppose the nonnegative orthant is translated by  $T \in \mathbb{R}^{m \times n}$ ; *id est*, consider  $\mathbb{R}_+^{m \times n} + T$ . Then projection on the translated orthant is [[122](#), §4.8]

$$H_{ij}^* = \max\{H_{ij}, T_{ij}\} \tag{2234}$$

□

**E.9.2.2.4 Example.** *Unique projection on truncated convex cone.*

Consider the problem of projecting a point  $x$  on a closed convex cone that is artificially bounded; really, a bounded convex polyhedron having a vertex at the origin:

$$\begin{aligned} &\underset{y \in \mathbb{R}^N}{\text{minimize}} \quad \|x - Ay\|_2 \\ &\text{subject to} \quad y \succeq 0 \\ &\quad \|y\|_\infty \leq 1 \end{aligned} \tag{2235}$$

where the convex cone has vertex-description ([§2.12.2.0.1](#)), for  $A \in \mathbb{R}^{n \times N}$

$$\mathcal{K} = \{Ay \mid y \succeq 0\} \tag{2236}$$

and where  $\|y\|_\infty \leq 1$  is the artificial bound. This is a convex optimization problem having no known closed-form solution, in general. It arises, for example, in the fitting of hearing aids designed around a programmable graphic equalizer (a filter bank whose only adjustable parameters are gain per frequency band each bounded above by unity). [[106](#)] [[107](#)] The problem is equivalent to a Schur-form semidefinite program ([§3.5.3](#))

$$\begin{aligned} &\underset{y \in \mathbb{R}^N, t \in \mathbb{R}}{\text{minimize}} \quad t \\ &\text{subject to} \quad \begin{bmatrix} tI & x - Ay \\ (x - Ay)^T & t \end{bmatrix} \succeq 0 \\ &\quad 0 \preceq y \preceq \mathbf{1} \end{aligned} \tag{2237}$$

□

### E.9.3 nonexpansivity

**E.9.3.0.1 Theorem.** *Nonexpansivity.*[[199](#), §2] [[122](#), §5.3]

When  $\mathcal{C} \subset \mathbb{R}^n$  is an arbitrary closed convex set, projector  $P$  projecting on  $\mathcal{C}$  is nonexpansive in the sense: for any vectors  $x, y \in \mathbb{R}^n$

$$\|Px - Py\| \leq \|x - y\| \tag{2238}$$

with equality when  $x - Px = y - Py$ . [E.19](#)

◊

---

[E.19](#)This condition for equality corrects an error in [[86](#)] (where the norm is applied to each side of the condition given here) easily revealed by counterexample.

**Proof.** [57]

$$\begin{aligned}\|x - y\|^2 &= \|Px - Py\|^2 + \|(I - P)x - (I - P)y\|^2 \\ &\quad + 2\langle x - Px, Px - Py \rangle + 2\langle y - Py, Py - Px \rangle\end{aligned}\tag{2239}$$

Nonnegativity of the last two terms follows directly from the *unique minimum-distance projection theorem* (§E.9.1.0.2).  $\diamond$

The foregoing proof reveals another flavor of nonexpansivity; for each and every  $x, y \in \mathbb{R}^n$

$$\|Px - Py\|^2 + \|(I - P)x - (I - P)y\|^2 \leq \|x - y\|^2\tag{2240}$$

Deutsch shows yet another: [122, §5.5]

$$\|Px - Py\|^2 \leq \langle x - y, Px - Py \rangle\tag{2241}$$

#### E.9.4 Easy projections

- To project any matrix  $H \in \mathbb{R}^{n \times n}$  orthogonally in Euclidean/Frobenius sense on subspace of symmetric matrices  $\mathbb{S}^n$  in isomorphic  $\mathbb{R}^{n^2}$ , take symmetric part of  $H$ ; (§2.2.2.0.1) *id est*,  $(H + H^T)/2$  is the projection.
- To project any  $H \in \mathbb{R}^{n \times n}$  orthogonally on symmetric hollow subspace  $\mathbb{S}_h^n$  in isomorphic  $\mathbb{R}^{n^2}$  (§2.2.3.0.1, §7.0.1), take symmetric part then zero all entries along main diagonal or *vice versa* (because this is projection on intersection of two subspaces); *id est*,  $(H + H^T)/2 - \delta^2(H)$ .
- To project a matrix on nonnegative orthant  $\mathbb{R}_+^{m \times n}$ , simply clip all negative entries to 0. Likewise, projection on nonpositive orthant  $\mathbb{R}_-^{m \times n}$  sees all positive entries clipped to 0. Projection on other orthants is equally simple with appropriate clipping.
- Projecting on hyperplane, halfspace, slab: §E.5.0.0.5.
- Projection of  $y \in \mathbb{R}^n$  on Euclidean ball  $\mathcal{B} = \{x \in \mathbb{R}^n \mid \|x - a\| \leq c\}$ : for  $y \neq a$ ,  $P_{\mathcal{B}}y = \frac{c}{\|y - a\|}(y - a) + a$ .
- Clipping in excess of  $|1|$ , each entry of point  $x \in \mathbb{R}^n$ , is equivalent to unique minimum-distance projection of  $x$  on a hypercube centered at the origin. (*confer* §E.10.3.2)
- Projection of  $x \in \mathbb{R}^n$  on a (rectangular) *hyperbox*: [65, §8.1.1]

$$\mathcal{C} = \{y \in \mathbb{R}^n \mid l \preceq y \preceq u, l \prec u\}\tag{2242}$$

$$P(x)_{k=0 \dots n} = \begin{cases} l_k, & x_k \leq l_k \\ x_k, & l_k \leq x_k \leq u_k \\ u_k, & x_k \geq u_k \end{cases}\tag{2243}$$

- Orthogonal projection of  $x$  on a *Cartesian subspace*, whose basis is some given subset of the Cartesian axes, zeroes entries corresponding to the remaining (complementary) axes.
- Projection of  $x$  on set of all cardinality- $k$  vectors  $\{y \mid \text{card } y \leq k\}$  keeps  $k$  entries of greatest magnitude and clips to 0 those remaining.

- Unique minimum-distance projection of  $H \in \mathbb{S}^n$  on positive semidefinite cone  $\mathbb{S}_+^n$ , in Euclidean/Frobenius sense, is accomplished by eigenvalue decomposition (diagonalization) followed by clipping all negative eigenvalues to 0.
- Unique minimum-distance projection on generally nonconvex subset of all matrices belonging to  $\mathbb{S}_+^n$  having rank not exceeding  $\rho$  (§2.9.2.1) is accomplished by clipping all negative eigenvalues to 0 and zeroing smallest nonnegative eigenvalues keeping only  $\rho$  largest. (§7.1.2)
- Unique minimum-distance projection, in Euclidean/Frobenius sense, of  $H \in \mathbb{R}^{m \times n}$  on the generally nonconvex subset of all  $m \times n$  matrices having rank no greater than  $k$  is the singular value decomposition (§A.6) of  $H$  having all singular values beyond its  $k^{\text{th}}$  zeroed. This is also a solution to projection in sense of spectral norm. [364, p.79, p.208]
- Projection of a real vector on the monotone nonnegative cone is identical to its projection on the monotone cone followed by clipping all negative entries of the result to 0. [307, §5]
- Projection on monotone nonnegative cone  $\mathcal{K}_{\mathcal{M}+} \subset \mathbb{R}^n$  in less than one cycle (in sense of alternating projections §E.10): [418].
- Fast projection on a simplicial cone: [425].
- Projection on closed convex cone  $\mathcal{K}$  of any point  $x \in -\mathcal{K}^*$ , belonging to polar cone, is equivalent to projection on origin. (§E.9.2)

$$\bullet P_{\mathbb{S}_+^N \cap \mathbb{S}_c^N} = P_{\mathbb{S}_+^N} P_{\mathbb{S}_c^N} \quad (\text{1410})$$

$$\bullet P_{\mathbb{R}_+^{N \times N} \cap \mathbb{S}_h^N} = P_{\mathbb{R}_+^{N \times N}} P_{\mathbb{S}_h^N} \quad (\text{§7.0.1.1})$$

$$\bullet P_{\mathbb{R}_+^{N \times N} \cap \mathbb{S}^N} = P_{\mathbb{R}_+^{N \times N}} P_{\mathbb{S}^N} \quad (\text{§E.9.5})$$

#### E.9.4.0.1 Exercise. Projection on spectral norm ball.

Find the unique minimum-distance projection on the convex set of all  $m \times n$  matrices whose largest singular value does not exceed 1; *id est*, on  $\{X \in \mathbb{R}^{m \times n} \mid \|X\|_2 \leq 1\}$  the spectral norm ball (§2.3.2.0.5). ▼

#### E.9.4.1 notes

Projection on Lorentz (second-order) cone: [65, exer.8.3c].

Deutsch [125] provides an algorithm for projection on polyhedral cones.

Youla [454, §2.5] lists eleven “useful projections”, of square-integrable uni- and bivariate real functions on various convex sets, in closed form.

Unique minimum-distance projection on an ellipsoid: Example 4.7.0.0.1.

Numerical algorithms for projection on the nonnegative simplex and 1-norm ball: [144].

Ferreira & Németh [158, cor.11] provide a formula for projection on a hyperplane in any vector norm.

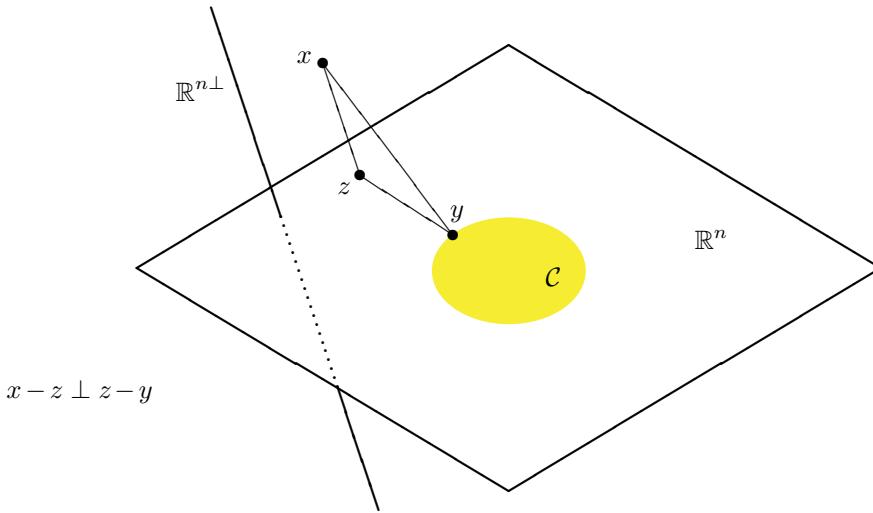


Figure 194: Closed convex set  $\mathcal{C}$  belongs to subspace  $\mathbb{R}^n$  (shown bounded in sketch and drawn without proper perspective). Point  $y$  is unique minimum-distance projection of  $x$  on  $\mathcal{C}$ ; equivalent to product of orthogonal projection of  $x$  on  $\mathbb{R}^n$  and minimum-distance projection of result  $z$  on  $\mathcal{C}$ .

### E.9.5 Projection on convex set in subspace

Suppose a convex set  $\mathcal{C}$  is wholly contained in some subspace  $\mathbb{R}^n$ . Then unique minimum-distance projection of any point in  $\mathbb{R}^n \oplus \mathbb{R}^{n\perp}$  on  $\mathcal{C}$  can be accomplished by first projecting orthogonally on that subspace, and then projecting the result on  $\mathcal{C}$ ; [122, §5.14] *id est*, the ordered product of two individual projections.

**Proof.** ( $\Leftarrow$ ) To show that, suppose unique minimum-distance projection  $P_{\mathcal{C}}x$  on  $\mathcal{C} \subset \mathbb{R}^n$  is  $y$  as illustrated in Figure 194;

$$\|x - y\| \leq \|x - q\| \quad \forall q \in \mathcal{C} \quad (2244)$$

Further suppose  $P_{\mathbb{R}^n}x$  equals  $z$ . By the *Pythagorean theorem*

$$\|x - y\|^2 = \|x - z\|^2 + \|z - y\|^2 \quad (2245)$$

because  $x - z \perp z - y$ . (2090) [280, §3.3] Then point  $y = P_{\mathcal{C}}x$  is the same as  $P_{\mathcal{C}}z$  because

$$\|z - y\|^2 = \|x - y\|^2 - \|x - z\|^2 \leq \|z - q\|^2 = \|x - q\|^2 - \|x - z\|^2 \quad \forall q \in \mathcal{C} \quad (2246)$$

which holds by assumption (2244).

( $\Rightarrow$ ) Now suppose  $z = P_{\mathbb{R}^n}x$  and

$$\|z - y\| \leq \|z - q\| \quad \forall q \in \mathcal{C} \quad (2247)$$

meaning  $y = P_{\mathcal{C}}z$ . Then point  $y$  is identical to  $P_{\mathcal{C}}x$  because

$$\|x - y\|^2 = \|x - z\|^2 + \|z - y\|^2 \leq \|x - q\|^2 = \|x - z\|^2 + \|z - q\|^2 \quad \forall q \in \mathcal{C} \quad (2248)$$

by assumption (2247). ♦

---

E.20 The goal, here, is to project on  $\mathcal{C}$  while remaining in that subspace.

This proof is extensible via translation argument. (§E.4) Unique minimum-distance projection on a convex set contained in an affine subset is, therefore, similarly accomplished.

Projecting matrix  $H \in \mathbb{R}^{n \times n}$  on convex cone  $\mathcal{K} = \mathbb{S}^n \cap \mathbb{R}_+^{n \times n}$  in isomorphic  $\mathbb{R}^{n^2}$  can be accomplished, for example, by first projecting on  $\mathbb{S}^n$  and only then projecting the result on  $\mathbb{R}_+^{n \times n}$  (confer §7.0.1). This is because projection product  $P_{\mathbb{R}_+^{n \times n}} P_{\mathbb{S}^n}$  is equivalent to projection on the subset of the nonnegative orthant in the symmetric matrix subspace.

## E.10 Alternating projection

Alternating projection is an iterative technique for finding a point in the intersection of a number of arbitrary closed convex sets  $\mathcal{C}_k$ , or for determining distance between two nonintersecting closed convex sets. Because it can sometimes be difficult or inefficient to compute the intersection or express it analytically, one naturally asks whether it is possible to instead project (unique minimum-distance) alternately on the individual  $\mathcal{C}_k$ ; often easier and what motivates adoption of this technique. Once a cycle of alternating projections (an *iteration*) is complete, we then *iterate* (repeat the cycle) until convergence. If the intersection of two closed convex sets is empty, then by *convergence* we mean the *iterate* (the result after a cycle of alternating projections) settles to a point of minimum distance separating the sets.

While alternating projection can find the point in the nonempty intersection closest to a given point  $b$ , it does not necessarily find the closest point. Finding that closest point is made dependable by an elegantly simple enhancement via correction to the alternating projection technique: this *Dykstra algorithm* (2286) for projection on the intersection is one of the most beautiful projection algorithms ever discovered. It is accurately interpreted as the discovery of what alternating projection originally sought to accomplish: unique minimum-distance projection on the nonempty intersection of a number of arbitrary closed convex sets  $\mathcal{C}_k$ . Alternating projection is, in fact, a special case of the Dykstra algorithm whose discussion we defer until §E.10.3.

### E.10.0.1 commutative projectors

A product of projection operators is generally not another projector. Given two arbitrary convex sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$  and their respective minimum-distance projection operators  $P_1$  and  $P_2$ : If projectors commute ( $P_1 P_2 = P_2 P_1$ ) for each and every  $x \in \mathbb{R}^n$ , then it is easy to show  $P_1 P_2 x \in \mathcal{C}_1 \cap \mathcal{C}_2$  and  $P_2 P_1 x \in \mathcal{C}_1 \cap \mathcal{C}_2$ . When projectors commute, their product is a projector; some point in the intersection can be found in a finite number of steps. While commutativity is a sufficient condition, it is not necessary; e.g., §6.8.1.1.

When  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are subspaces, in particular, projectors  $P_1$  and  $P_2$  commute if and only if  $P_1 P_2 = P_{\mathcal{C}_1 \cap \mathcal{C}_2}$  or iff  $P_2 P_1 = P_{\mathcal{C}_1 \cap \mathcal{C}_2}$  or iff  $P_1 P_2$  is the orthogonal projection on a Euclidean subspace. [122, lem.9.2] Subspace projectors will commute, for example, when  $P_1(\mathcal{C}_2) \subset \mathcal{C}_2$  or  $P_2(\mathcal{C}_1) \subset \mathcal{C}_1$  or  $\mathcal{C}_1 \subset \mathcal{C}_2$  or  $\mathcal{C}_2 \subset \mathcal{C}_1$  or  $\mathcal{C}_1 \perp \mathcal{C}_2$ . When subspace projectors commute, this means we can find a point from the intersection of those subspaces in a finite number of steps; in fact, the closest point. Orthogonal projection on orthogonal subspaces (or intersecting orthogonal affine subsets), in particular, can be performed in any order to find the closest point in their intersection in a number of steps = number of subspaces (or affine subsets).

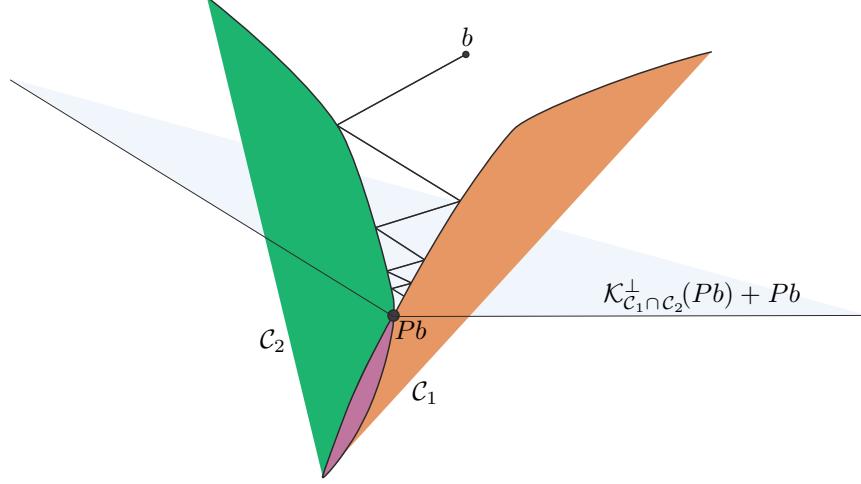


Figure 195: First several alternating projections, in von Neumann-style projection (2259) of point  $b$ , converging on closest point  $Pb$  • in intersection of two closed convex sets in  $\mathbb{R}^2$ :  $C_1$  and  $C_2$  partially drawn in vicinity of their intersection. For this particular example, it is possible to start anywhere in a large neighborhood of  $b$  and still converge to  $Pb$ . Pointed normal cone  $K^\perp$  at  $Pb$  (454) is translated to  $Pb$  the unique minimum-distance projection of  $b$  on intersection. Alternating projections are themselves robust with respect to significant noise because they belong to this translated normal cone.

#### E.10.0.2 noncommutative projectors

Typically, one considers the method of alternating projection when projectors do not commute; *id est*, when  $P_1P_2 \neq P_2P_1$ .

The iconic example for noncommutative projectors illustrated in Figure 195 shows the iterates converging to the closest point in the intersection of two arbitrary convex sets. Yet simple examples like Figure 196 reveal that noncommutative alternating projection does not always yield the closest point, although we shall show it always yields some point in the intersection or a point that attains the distance between two convex sets.

Alternating projection is also known as *successive projection* [202] [199] [67], *cyclic projection* [168] [291, §3.2], *successive approximation* [86], or *projection on convex sets* [361] [362, §6.4]. It is traced back to von Neumann, 1933 [409], and later Wiener [438] who showed that higher iterates of a product of two orthogonal projections on subspaces converge at each point in the ambient space to the unique minimum-distance projection on the intersection of the two subspaces. More precisely, if  $\mathcal{R}_1$  and  $\mathcal{R}_2$  are closed subspaces of a Euclidean space and  $P_1$  and  $P_2$  respectively denote orthogonal projection on  $\mathcal{R}_1$  and  $\mathcal{R}_2$ , then for each vector  $b$  in that space,

$$\lim_{i \rightarrow \infty} (P_1P_2)^i b = P_{\mathcal{R}_1 \cap \mathcal{R}_2} b \quad (2249)$$

Deutsch [122, thm.9.8, thm.9.35] shows rate of convergence for subspaces to be *geometric* [453, §1.4.4]; bounded above by  $\kappa^{2i+1}\|b\|$ ,  $i=0, 1, 2, \dots$ , where  $0 \leq \kappa < 1$ :

$$\|(P_1P_2)^i b - P_{\mathcal{R}_1 \cap \mathcal{R}_2} b\| \leq \kappa^{2i+1}\|b\| \quad (2250)$$

This means convergence can be slow when  $\kappa$  is close to 1. Rate of convergence on intersecting halfspaces is also geometric. [123] [331]

This von Neumann sense of alternating projection may be applied to convex sets that are not subspaces, although convergence is not necessarily to the unique minimum-distance

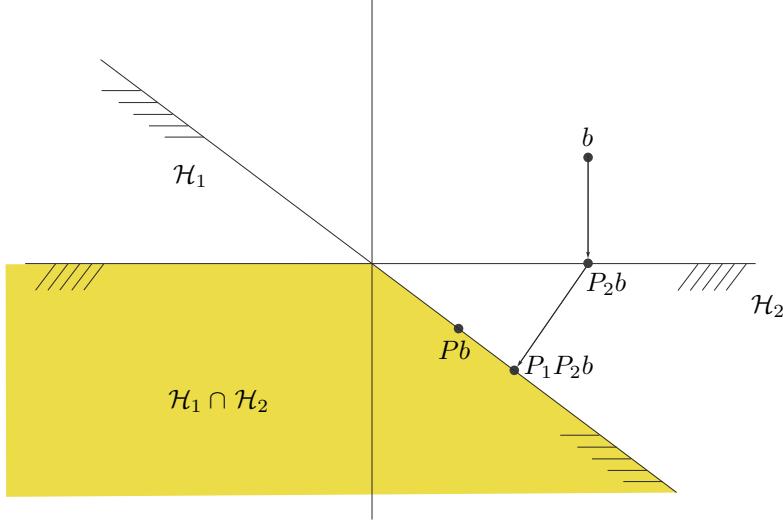


Figure 196: The sets  $\{\mathcal{C}_k\}$  in this example comprise two halfspaces  $\mathcal{H}_1$  and  $\mathcal{H}_2$ . The von Neumann-style alternating projection in  $\mathbb{R}^2$  quickly converges to  $P_1P_2b$  (feasibility). Unique minimum-distance projection on the intersection is, of course,  $Pb$ .

projection on the intersection. Figure 195 illustrates one application where convergence is reasonably geometric and the result is the unique minimum-distance projection. Figure 196, in contrast, demonstrates convergence in one iteration to a fixed point (of the projection product)<sup>E.21</sup> in the intersection of two halfspaces; a.k.a, feasibility problem. It was Dykstra who in 1983 [145] (§E.10.3) first solved this projection problem.

#### E.10.0.3 the bullets

Alternating projection has, therefore, various meaning dependent on the application or field of study; it may be interpreted to be: a distance problem, a feasibility problem (von Neumann), or a projection problem (Dykstra):

- **Distance.** Figure 197a-b. Find a unique point of projection  $P_1b \in \mathcal{C}_1$  that attains the distance between any two closed convex sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$ ;

$$\|P_1b - b\| = \text{dist}(\mathcal{C}_1, \mathcal{C}_2) \triangleq \inf_{z \in \mathcal{C}_2} \|P_1z - z\| \quad (2251)$$

- **Feasibility.** Figure 197c,  $\bigcap \mathcal{C}_k \neq \emptyset$ . Given a number  $L$  of indexed closed convex sets  $\mathcal{C}_k \subset \mathbb{R}^n$ , find any fixed point in their intersection by iterating (i) a projection product starting from  $b$ ;

$$\left( \prod_{i=1}^{\infty} \prod_{k=1}^L P_k \right) b \in \bigcap_{k=1}^L \mathcal{C}_k \quad (2252)$$

- **Optimization.** Figure 197c,  $\bigcap \mathcal{C}_k \neq \emptyset$ . Given a number of indexed closed convex sets  $\mathcal{C}_k \subset \mathbb{R}^n$ , uniquely project a given point  $b$  on  $\bigcap \mathcal{C}_k$ ;

$$\|Pb - b\| = \inf_{x \in \bigcap \mathcal{C}_k} \|x - b\| \quad (2253)$$

---

<sup>E.21</sup>Fixed point of mapping  $T: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a point  $x$  whose image is identical under the map; id est,  $Tx = x$ .

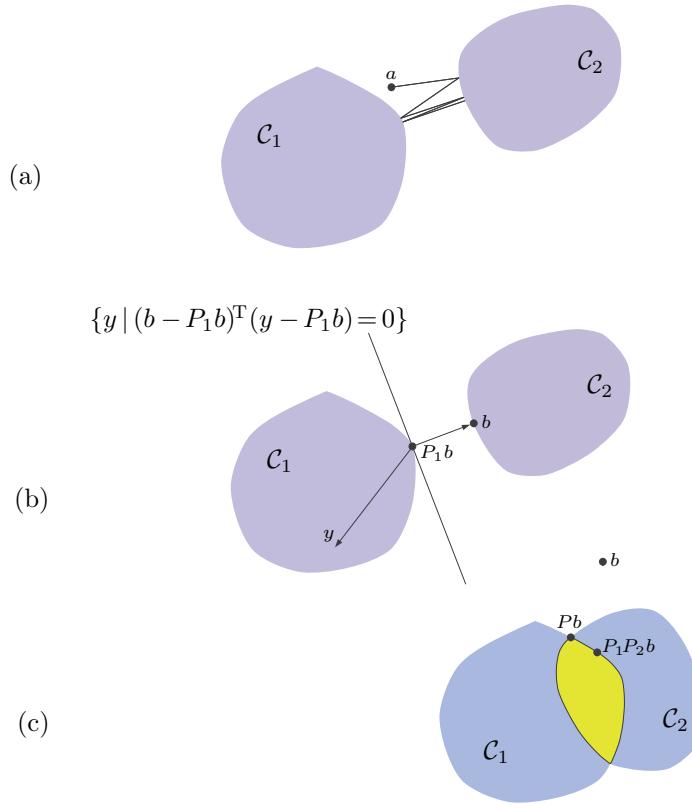


Figure 197:

**(a) (distance)** Intersection of two convex sets in  $\mathbb{R}^2$  is empty. Method of alternating projection would be applied to find that point in  $\mathcal{C}_1$  nearest  $\mathcal{C}_2$ .

**(b) (distance)** Given  $b \in \mathcal{C}_2$ , by *projection theorem E.9.1.0.2*, point  $\bullet P_1 b \in \mathcal{C}_1$  is nearest  $b$  iff  $(b - P_1 b)^T(y - P_1 b) \leq 0 \forall y \in \mathcal{C}_1$ . When  $P_1 b$  attains distance between the two sets, hyperplane  $\{y \mid (b - P_1 b)^T(y - P_1 b) = 0\}$  separates  $\mathcal{C}_1$  from  $\mathcal{C}_2$ . [65, §2.5.1]

**(c) (0 distance)** Intersection is nonempty; distance between sets equals 0.

**(feasibility)** We may just want a fixed point of projection product  $P_1 P_2 b$  in  $\bigcap \mathcal{C}_k$ .

**(optimization)** Or we may want that point  $Pb$  in  $\bigcap \mathcal{C}_k$  nearest point  $b$ .

### E.10.1 Distance and existence

Existence of a fixed point is established:

#### E.10.1.0.1 Theorem. *Distance.*

[86]

Given any two closed convex sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$  in  $\mathbb{R}^n$ , then  $P_1 b \in \mathcal{C}_1$  is a fixed point of projection product  $P_1 P_2$  if and only if  $P_1 b$  is a point of  $\mathcal{C}_1$  nearest  $\mathcal{C}_2$ .  $\diamond$

**Proof.** ( $\Rightarrow$ ) Given fixed point  $a = P_1 P_2 a \in \mathcal{C}_1$  with  $b \triangleq P_2 a \in \mathcal{C}_2$  in tandem so that  $a = P_1 b$ , then by the *unique minimum-distance projection theorem* (§E.9.1.0.2)

$$\begin{aligned} (b - a)^T(u - a) &\leq 0 \quad \forall u \in \mathcal{C}_1 \\ (a - b)^T(v - b) &\leq 0 \quad \forall v \in \mathcal{C}_2 \\ &\Leftrightarrow \\ \|a - b\| &\leq \|u - v\| \quad \forall u \in \mathcal{C}_1 \text{ and } \forall v \in \mathcal{C}_2 \end{aligned} \quad (2254)$$

by Cauchy-Schwarz inequality [343]

$$|\langle x, y \rangle| \leq \|x\| \|y\| \quad (2255)$$

(with equality iff  $x = \kappa y$  where  $\kappa \in \mathbb{R}$  (33) [254, p.137]).

( $\Leftarrow$ ) Suppose  $a \in \mathcal{C}_1$  and  $\|a - P_2 a\| \leq \|u - P_2 u\| \forall u \in \mathcal{C}_1$  and we choose  $u = P_1 P_2 a$ . Then

$$\|u - P_2 u\| = \|P_1 P_2 a - P_2 P_1 P_2 a\| \leq \|a - P_2 a\| \Leftrightarrow a = P_1 P_2 a \quad (2256)$$

Thus  $a = P_1 b$  (with  $b = P_2 a \in \mathcal{C}_2$ ) is a fixed point in  $\mathcal{C}_1$  of the projection product  $P_1 P_2$ .  $\blacklozenge$

### E.10.2 Feasibility and convergence

The set of all fixed points of any nonexpansive mapping is a closed convex set. [177, lem.3.4] [29, §1] The projection product  $P_1 P_2$  is nonexpansive by Theorem E.9.3.0.1 because, for any vectors  $x, a \in \mathbb{R}^n$

$$\|P_1 P_2 x - P_1 P_2 a\| \leq \|P_2 x - P_2 a\| \leq \|x - a\| \quad (2257)$$

If the intersection of two closed convex sets  $\mathcal{C}_1 \cap \mathcal{C}_2$  is empty, then the iterates converge to a point of minimum distance, a fixed point of the projection product. Otherwise, convergence is to some fixed point in their intersection (a feasible solution) whose existence is guaranteed by virtue of the fact that each and every point in the convex intersection is in one-to-one correspondence with fixed points of the nonexpansive projection product.

Bauschke & Borwein [29, §2] argue that any sequence monotonic in the sense of Fejér is convergent:  $\blacklozenge$

#### E.10.2.0.1 Definition. *Fejér monotonicity.*

[301]

Given closed convex set  $\mathcal{C} \neq \emptyset$ , then a sequence  $x_i \in \mathbb{R}^n$ ,  $i = 0, 1, 2, \dots$ , is monotonic in the sense of Fejér with respect to  $\mathcal{C}$  iff

$$\|x_{i+1} - c\| \leq \|x_i - c\| \quad \text{for all } i \geq 0 \text{ and each and every } c \in \mathcal{C} \quad (2258)$$

 $\triangle$ 

**E.22** Point  $b = P_2 a$  can be shown, similarly, to be a fixed point of product  $P_2 P_1$ .

**E.23** Other authors prove convergence by different means; e.g., [199] [67].

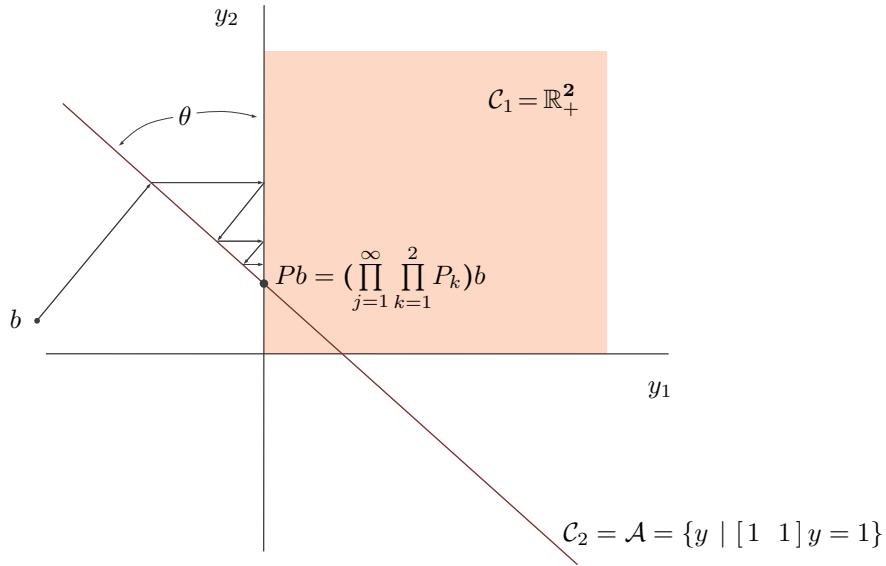


Figure 198: From Example E.10.2.0.2 in  $\mathbb{R}^2$  showing von Neumann-style alternating projection to find feasible solution belonging to intersection of nonnegative orthant with hyperplane  $\mathcal{A}$ . Point  $Pb$  lies at intersection of hyperplane with ordinate axis. In this particular example, feasible solution found is coincidentally optimal. Rate of convergence depends upon angle  $\theta$ ; as it becomes more acute, convergence slows. [199, §3]

Given  $x_0 \triangleq b$ , if we express each iteration of alternating projection by

$$x_{i+1} = P_1 P_2 x_i , \quad i = 0, 1, 2 \dots \quad (2259)$$

and define any fixed point  $a = P_1 P_2 a$ , then sequence  $x_i$  is Fejér monotone with respect to fixed point  $a$  because

$$\|P_1 P_2 x_i - a\| \leq \|x_i - a\| \quad \forall i \geq 0 \quad (2260)$$

by nonexpansivity. The nonincreasing sequence  $\|P_1 P_2 x_i - a\|$  is bounded below hence convergent because any bounded monotonic sequence in  $\mathbb{R}$  is convergent; [289, §1.2] [43, §1.1]  $P_1 P_2 x_{i+1} = P_1 P_2 x_i = x_{i+1}$ . Sequence  $x_i$  therefore converges to some fixed point. If the intersection  $\mathcal{C}_1 \cap \mathcal{C}_2$  is nonempty, convergence is to some point there by the *distance theorem*. Otherwise,  $x_i$  converges to a point in  $\mathcal{C}_1$  of minimum distance to  $\mathcal{C}_2$ .

#### E.10.2.0.2 Example. Hyperplane/orthant intersection.

Find a feasible solution (2252) belonging to the nonempty intersection of two convex sets: given full-rank  $A \in \mathbb{R}^{m \times n}$  and  $\beta \in \mathcal{R}(A)$

$$\mathcal{C}_1 \cap \mathcal{C}_2 = \mathbb{R}_+^n \cap \mathcal{A} = \{y \mid y \succeq 0\} \cap \{y \mid Ay = \beta\} \subset \mathbb{R}^n \quad (2261)$$

the nonnegative orthant intersecting affine subset  $\mathcal{A}$  (an intersection of hyperplanes). Projection of an iterate  $x_i \in \mathbb{R}^n$  on  $\mathcal{A}$  is calculated

$$P_2 x_i = x_i - A^T (AA^T)^{-1} (Ax_i - \beta) \quad (2143)$$

while, thereafter, projection of the result on the orthant is simply

$$x_{i+1} = P_1 P_2 x_i = \max\{\mathbf{0}, P_2 x_i\} \quad (2262)$$

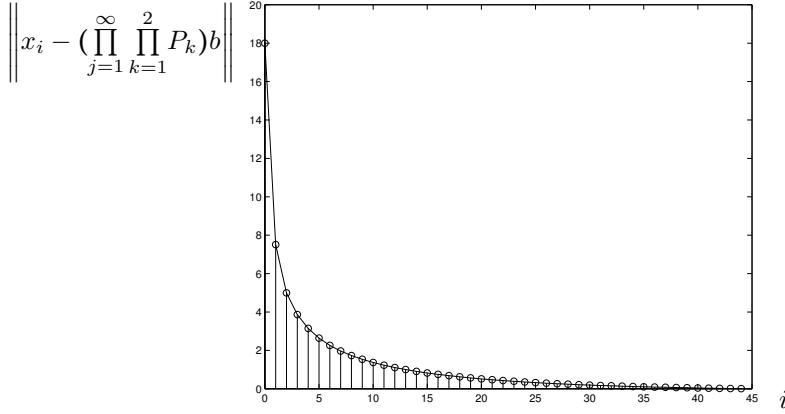


Figure 199: Example E.10.2.0.2 in  $\mathbb{R}^{1000}$ ; geometric convergence of iterates in norm.

where the maximum is entrywise (§E.9.2.2.3).

One realization of this problem in  $\mathbb{R}^2$  is illustrated in Figure 198: For  $A = [1 \ 1]$ ,  $\beta = 1$ , and  $x_0 = b = [-3 \ 1/2]^T$ , iterates converge to a feasible solution  $Pb = [0 \ 1]^T$ .

To give a more palpable sense of convergence in higher dimension, we do this example again but now we compute an alternating projection for the case  $A \in \mathbb{R}^{400 \times 1000}$ ,  $\beta \in \mathbb{R}^{400}$ , and  $b \in \mathbb{R}^{1000}$ , all of whose entries are independently and randomly set to a uniformly distributed real number in the interval  $[-1, 1]$ . Convergence is illustrated in Figure 199.  $\square$

This application of alternating projection to feasibility is extensible to any finite number of closed convex sets.

#### E.10.2.0.3 Example. Under- and over-projection.

[62, §3]

Consider the following variation of alternating projection: We begin with some point  $x_0 \in \mathbb{R}^n$  then project that point on convex set  $\mathcal{C}$  and then project that same point  $x_0$  on convex set  $\mathcal{D}$ . To the first iterate we assign  $x_1 = \frac{1}{2}(P_{\mathcal{C}}(x_0) + P_{\mathcal{D}}(x_0))$ . More generally,

$$x_{i+1} = \frac{1}{2}(P_{\mathcal{C}}(x_i) + P_{\mathcal{D}}(x_i)), \quad i = 0, 1, 2, \dots \quad (2263)$$

Because the Cartesian product of convex sets remains convex, (§2.1.8) we can reformulate this problem.

Consider the convex set

$$\mathcal{Z} \triangleq \begin{bmatrix} \mathcal{C} \\ \mathcal{D} \end{bmatrix} \quad (2264)$$

representing Cartesian product  $\mathcal{C} \times \mathcal{D}$ . Now, those two projections  $P_{\mathcal{C}}$  and  $P_{\mathcal{D}}$  are equivalent to one projection on the Cartesian product; *id est*,

$$P_{\mathcal{Z}} \left( \begin{bmatrix} x_i \\ x_i \end{bmatrix} \right) = \begin{bmatrix} P_{\mathcal{C}}(x_i) \\ P_{\mathcal{D}}(x_i) \end{bmatrix} \quad (2265)$$

Define the subspace

$$\mathcal{R} \triangleq \left\{ v \in \begin{bmatrix} \mathbb{R}^n \\ \mathbb{R}^n \end{bmatrix} \mid [I \ -I]v = \mathbf{0} \right\} \quad (2266)$$

By the results in Example E.5.0.0.7

$$P_{\mathcal{R}\mathcal{Z}} \left( \begin{bmatrix} x_i \\ x_i \end{bmatrix} \right) = P_{\mathcal{R}} \left( \begin{bmatrix} P_{\mathcal{C}}(x_i) \\ P_{\mathcal{D}}(x_i) \end{bmatrix} \right) = \frac{1}{2} \begin{bmatrix} P_{\mathcal{C}}(x_i) + P_{\mathcal{D}}(x_i) \\ P_{\mathcal{C}}(x_i) + P_{\mathcal{D}}(x_i) \end{bmatrix} \quad (2267)$$

This means the proposed variation of alternating projection is equivalent to an alternation of projection on convex sets  $\mathcal{Z}$  and  $\mathcal{R}$ . If  $\mathcal{Z}$  and  $\mathcal{R}$  intersect, these iterations will converge to a point in their intersection; hence, to a point in the intersection of  $\mathcal{C}$  and  $\mathcal{D}$ .

We need not apply equal weighting to the projections, as supposed in (2263). In that case, definition of  $\mathcal{R}$  would change accordingly.  $\square$

### E.10.2.1 Relative measure of convergence

Inspired by Fejér monotonicity, the alternating projection algorithm from the example of convergence illustrated by Figure 199 employs a redundant sequence: The first sequence (indexed by  $j$ ) estimates point  $(\prod_{j=1}^{\infty} \prod_{k=1}^L P_k)b$  in the presumably nonempty intersection of  $L$  convex sets, then the quantity

$$\left\| x_i - \left( \prod_{j=1}^{\infty} \prod_{k=1}^L P_k \right) b \right\| \quad (2268)$$

in second sequence  $x_i$  is observed per iteration  $i$  for convergence. *A priori* knowledge of a feasible solution (2252) is both impractical and antithetical. We need another measure:

Nonexpansivity implies

$$\left\| \left( \prod_{\ell=1}^L P_{\ell} \right) x_{k,i-1} - \left( \prod_{\ell=1}^L P_{\ell} \right) x_{ki} \right\| = \|x_{ki} - x_{k,i+1}\| \leq \|x_{k,i-1} - x_{ki}\| \quad (2269)$$

where

$$x_{ki} \triangleq P_k x_{k+1,i} \in \mathbb{R}^n, \quad x_{L+1,i} \triangleq x_{1,i-1} \quad (2270)$$

$x_{ki}$  represents unique minimum-distance projection of  $x_{k+1,i}$  on convex set  $k$  at iteration  $i$ . So a good convergence measure is total monotonic sequence

$$\varepsilon_i \triangleq \sum_{k=1}^L \|x_{ki} - x_{k,i+1}\|, \quad i=0,1,2\dots \quad (2271)$$

where  $\lim_{i \rightarrow \infty} \varepsilon_i = 0$  whether or not the intersection is nonempty.

**E.10.2.1.1 Example.** *Affine subset  $\cap$  positive semidefinite cone:  $\mathcal{A} \cap \mathbb{S}_+^n$ .*  
Consider the problem of finding  $X \in \mathbb{S}^n$  that satisfies

$$X \succeq 0, \quad \langle A_j, X \rangle = b_j, \quad j=1\dots m \quad (2272)$$

given nonzero  $A_j \in \mathbb{S}^n$  and real  $b_j$ . Here we take  $\mathcal{C}_1$  to be the positive semidefinite cone  $\mathbb{S}_+^n$  while  $\mathcal{C}_2$  is the affine subset of  $\mathbb{S}^n$

$$\begin{aligned} \mathcal{C}_2 &= \mathcal{A} \triangleq \{X \mid \langle A_j, X \rangle = b_j, j=1\dots m\} \subseteq \mathbb{S}^n \\ &= \{X \mid \text{tr}(A_j X) = b_j, j=1\dots m\} \\ &= \{X \mid \begin{bmatrix} \text{svec}(A_1)^T \\ \vdots \\ \text{svec}(A_m)^T \end{bmatrix} \text{svec } X = b\} \\ &\triangleq \{X \mid A \text{ svec } X = b\} \end{aligned} \quad (2273)$$

where  $b = [b_j] \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{m \times n(n+1)/2}$ , and symmetric vectorization svec is defined by (57). Projection of iterate  $X_i \in \mathbb{S}^n$  on  $\mathcal{A}$  is: (§E.5.0.0.7)

$$P_2 \text{svec } X_i = \text{svec } X_i - A^\dagger (A \text{svec } X_i - b) \quad (2274)$$

Euclidean distance from  $X_i$  to  $\mathcal{A}$  is therefore

$$\text{dist}(X_i, \mathcal{A}) = \|X_i - P_2 X_i\|_{\text{F}} = \|A^\dagger (A \text{svec } X_i - b)\|_2 \quad (2275)$$

Projection of  $P_2 X_i \triangleq \sum_j \lambda_j q_j q_j^T$  on the positive semidefinite cone (§7.1.2) is found from its eigenvalue decomposition (§A.5.1);

$$P_1 P_2 X_i = \sum_{j=1}^n \max\{0, \lambda_j\} q_j q_j^T \quad (2276)$$

Distance from  $P_2 X_i$  to the positive semidefinite cone is therefore

$$\text{dist}(P_2 X_i, \mathbb{S}_+^n) = \|P_2 X_i - P_1 P_2 X_i\|_{\text{F}} = \sqrt{\sum_{j=1}^n (\min\{0, \lambda_j\})^2} \quad (2277)$$

When the intersection is empty,  $\mathcal{A} \cap \mathbb{S}_+^n = \emptyset$ , the iterates converge to that positive semidefinite matrix closest to  $\mathcal{A}$  in the Euclidean sense. Otherwise, convergence is to some point in the nonempty intersection.

Barvinok (§2.9.3.0.1) shows that if a solution to (2272) exists, then there exists an  $X \in \mathcal{A} \cap \mathbb{S}_+^n$  such that

$$\text{rank } X \leq \left\lfloor \frac{\sqrt{8m+1}-1}{2} \right\rfloor \quad (275)$$

□

### E.10.2.1.2 Example. Semidefinite matrix completion.

Continuing Example E.10.2.1.1: When  $m \leq n(n+1)/2$  and the  $A_j$  matrices are distinct members of the standard orthonormal basis  $\{E_{\ell q} \in \mathbb{S}^n\}$  (60)

$$\{A_j \in \mathbb{S}^n, j=1 \dots m\} \subseteq \{E_{\ell q}\} = \left\{ \begin{array}{ll} e_\ell e_\ell^T, & \ell = q = 1 \dots n \\ \frac{1}{\sqrt{2}}(e_\ell e_q^T + e_q e_\ell^T), & 1 \leq \ell < q \leq n \end{array} \right\} \quad (2278)$$

and when the constants  $b_j$  are set to constrained entries of variable  $X \triangleq [X_{\ell q}] \in \mathbb{S}^n$

$$\{b_j, j=1 \dots m\} \subseteq \left\{ \begin{array}{ll} X_{\ell q}, & \ell = q = 1 \dots n \\ X_{\ell q} \sqrt{2}, & 1 \leq \ell < q \leq n \end{array} \right\} = \{\langle X, E_{\ell q} \rangle\} \quad (2279)$$

then the equality constraints in (2272) fix individual entries of  $X \in \mathbb{S}^n$ . Thus the feasibility problem becomes a *positive semidefinite matrix completion problem*. Projection of iterate  $X_i \in \mathbb{S}^n$  on  $\mathcal{A}$  simplifies to (confer (2274))

$$P_2 \text{svec } X_i = \text{svec } X_i - A^T (A \text{svec } X_i - b) \quad (2280)$$

From this we can see that orthogonal projection is achieved simply by setting corresponding entries of  $P_2 X_i$  to the known entries of  $X$ , while the entries of  $P_2 X_i$  remaining are set to corresponding entries of the current iterate  $X_i$ .

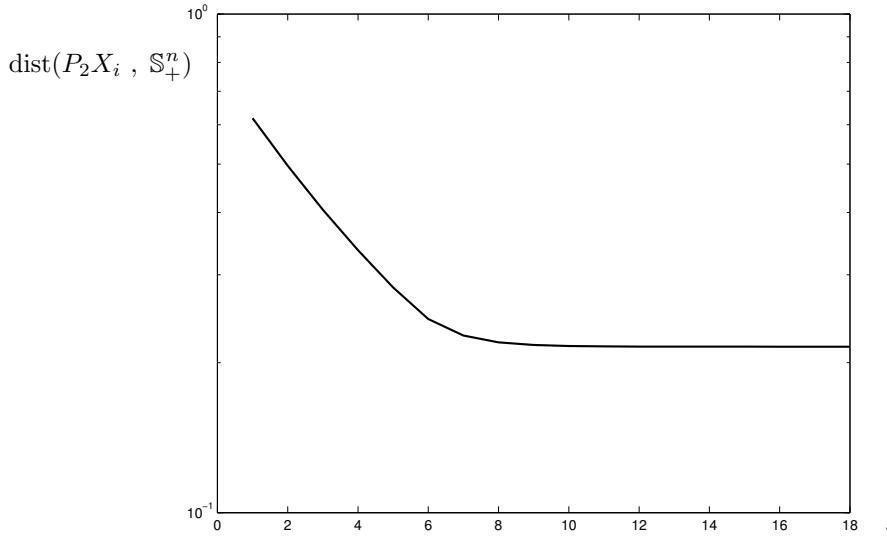


Figure 200: Distance ([confer\(2277\)](#)) between PSD cone and iterate ([2280](#)) in affine subset  $\mathcal{A}$  ([2273](#)) for Laurent's completion problem; initially, decreasing geometrically.

Using this technique, we find a positive semidefinite completion for

$$\begin{bmatrix} 4 & 3 & ? & 2 \\ 3 & 4 & 3 & ? \\ ? & 3 & 4 & 3 \\ 2 & ? & 3 & 4 \end{bmatrix} \quad (2281)$$

Initializing the unknown entries to 0, they all converge geometrically to 1.5858 (rounded) after about 42 iterations.

Laurent gives a problem for which no positive semidefinite completion exists: [\[265\]](#)

$$\begin{bmatrix} 1 & 1 & ? & 0 \\ 1 & 1 & 1 & ? \\ ? & 1 & 1 & 1 \\ 0 & ? & 1 & 1 \end{bmatrix} \quad (2282)$$

Initializing unknowns to 0, by alternating projection we find the constrained matrix closest to the positive semidefinite cone,

$$\begin{bmatrix} 1 & 1 & 0.5454 & 0 \\ 1 & 1 & 1 & 0.5454 \\ 0.5454 & 1 & 1 & 1 \\ 0 & 0.5454 & 1 & 1 \end{bmatrix} \quad (2283)$$

and we find the positive semidefinite matrix closest to affine subset  $\mathcal{A}$  ([2273](#)):

$$\begin{bmatrix} 1.0521 & 0.9409 & 0.5454 & 0.0292 \\ 0.9409 & 1.0980 & 0.9451 & 0.5454 \\ 0.5454 & 0.9451 & 1.0980 & 0.9409 \\ 0.0292 & 0.5454 & 0.9409 & 1.0521 \end{bmatrix} \quad (2284)$$

These matrices ([2283](#)) and ([2284](#)) attain the Euclidean distance  $\text{dist}(\mathcal{A}, \mathbb{S}_+^n)$ . Convergence is illustrated in Figure 200.  $\square$

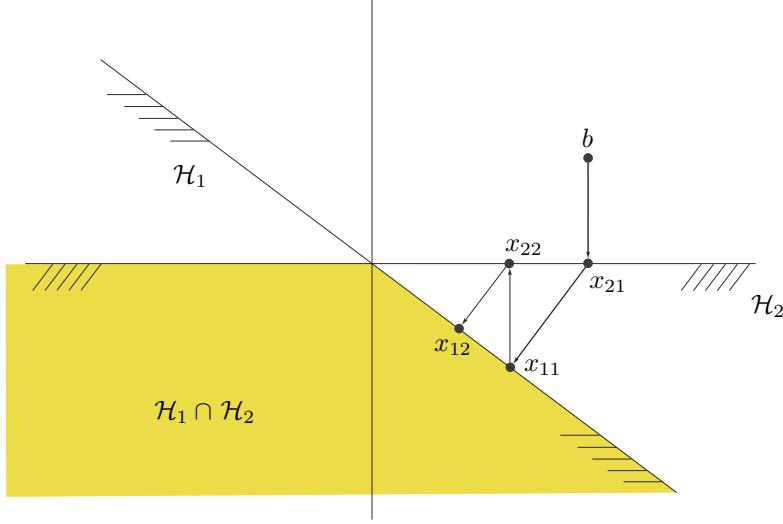


Figure 201:  $\mathcal{H}_1$  and  $\mathcal{H}_2$  are the same halfspaces as in Figure 196. Dykstra's alternating projection algorithm generates the alternations  $b, x_{21}, x_{11}, x_{22}, x_{12}, x_{21}, \dots, x_{12}$ . The path illustrated from  $b$  to  $x_{12}$  in  $\mathbb{R}^2$  terminates at the desired result:  $Pb$  in Figure 196. The  $\{y_{ki}\}$  correspond to the first two difference vectors drawn (in the first iteration  $i=1$ ), then oscillate between zero and a negative vector thereafter. These alternations are not so robust in presence of noise as for the example in Figure 195.

### E.10.3 Optimization and projection

Unique projection on the nonempty intersection of arbitrary convex sets, to find the closest point therein, is a convex optimization problem. The first successful application of alternating projection to this problem is attributed to Dykstra [145] [66] who in 1983 provided an elegant algorithm that prevails today. In 1988, Han [202] rediscovered the algorithm and provided a primal–dual convergence proof. A synopsis of the history of alternating projection<sup>E.24</sup> can be found in [68] where it becomes apparent that Dykstra's work is seminal; his algorithm appears in work as diverse as *machine control*. [205, §5]

#### E.10.3.1 Dykstra's algorithm

Assume we are given some point  $b \in \mathbb{R}^n$  and closed convex sets  $\{\mathcal{C}_k \subset \mathbb{R}^n \mid k=1 \dots L\}$ . Let  $x_{ki} \in \mathbb{R}^n$  and  $y_{ki} \in \mathbb{R}^n$  respectively denote a *primal* and *dual vector* (whose meaning can be deduced from Figure 201 and Figure 202) associated with set  $k$  at iteration  $i$ . Initialize

$$y_{k0} = 0 \quad \forall k=1 \dots L \quad \text{and} \quad x_{1,0} = b \tag{2285}$$

<sup>E.24</sup>For a synopsis of alternating projection applied to distance geometry, see [394, §3.1].

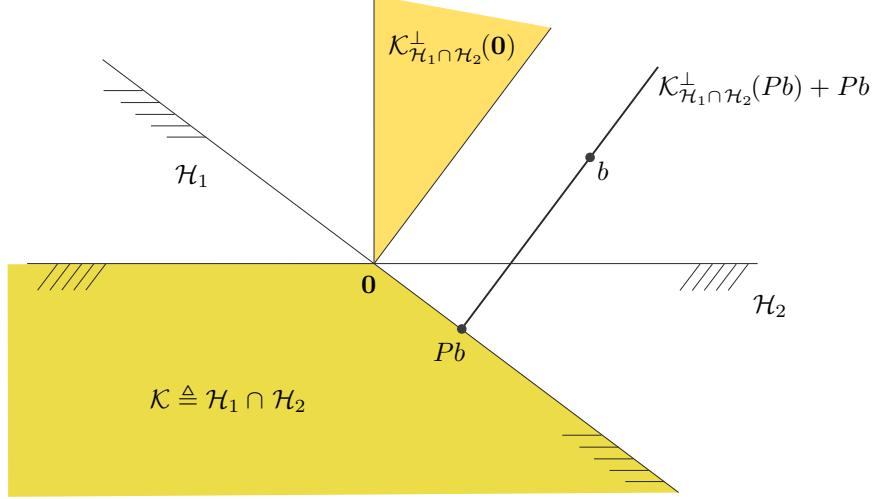


Figure 202: Two examples (truncated): normal cone to  $\mathcal{H}_1 \cap \mathcal{H}_2$  at the origin and at boundary point  $Pb$ .  $\mathcal{H}_1$  and  $\mathcal{H}_2$  are the same halfspaces from Figure 201. Normal cone at origin  $\mathcal{K}_{\mathcal{H}_1 \cap \mathcal{H}_2}^\perp(\mathbf{0})$  is simply  $-\mathcal{K}^*$ .

Denoting by  $P_k t$  the unique minimum-distance projection of  $t$  on  $\mathcal{C}_k$ , and for convenience  $x_{L+1,i} = x_{1,i-1}$  (2270), iterate  $x_{1i}$  calculation proceeds:<sup>E.25</sup>

$$\begin{aligned}
 & \text{for } i=1, 2, \dots \text{ until convergence } \{ \\
 & \quad \text{for } k=L \dots 1 \{ \\
 & \quad \quad t = x_{k+1,i} - y_{k,i-1} \\
 & \quad \quad x_{ki} = P_k t \\
 & \quad \quad y_{ki} = P_k t - t \\
 & \quad \}
 \end{aligned} \tag{2286}$$

Assuming a nonempty intersection, then iterates converge to the unique minimum-distance projection of point  $b$  on that intersection; [122, §9.24]

$$Pb = \lim_{i \rightarrow \infty} x_{1i} \tag{2287}$$

In the case that all  $\mathcal{C}_k$  are affine, then calculation of  $y_{ki}$  is superfluous and the algorithm becomes identical to alternating projection. [122, §9.26] [168, §1] Dykstra's algorithm is so simple, elegant, and represents such a tiny increment in computational intensity over alternating projection, it is nearly always arguably cost effective.

### E.10.3.2 Normal cone

Glunt [176, §4] observes that the overall effect of Dykstra's iterative procedure is to drive  $t$  toward the translated normal cone to  $\bigcap \mathcal{C}_k$  at the solution  $Pb$  (translated to  $Pb$ ). The normal cone gets its name from its graphical construction; which is, loosely speaking, to draw the outward-normals at  $Pb$  (Definition E.9.1.0.1) to all the convex sets  $\mathcal{C}_k$  touching  $Pb$ . Relative interior of the normal cone subtends these normal vectors.

<sup>E.25</sup>We reverse order of projection ( $k=L \dots 1$ ) in the algorithm for continuity of exposition.

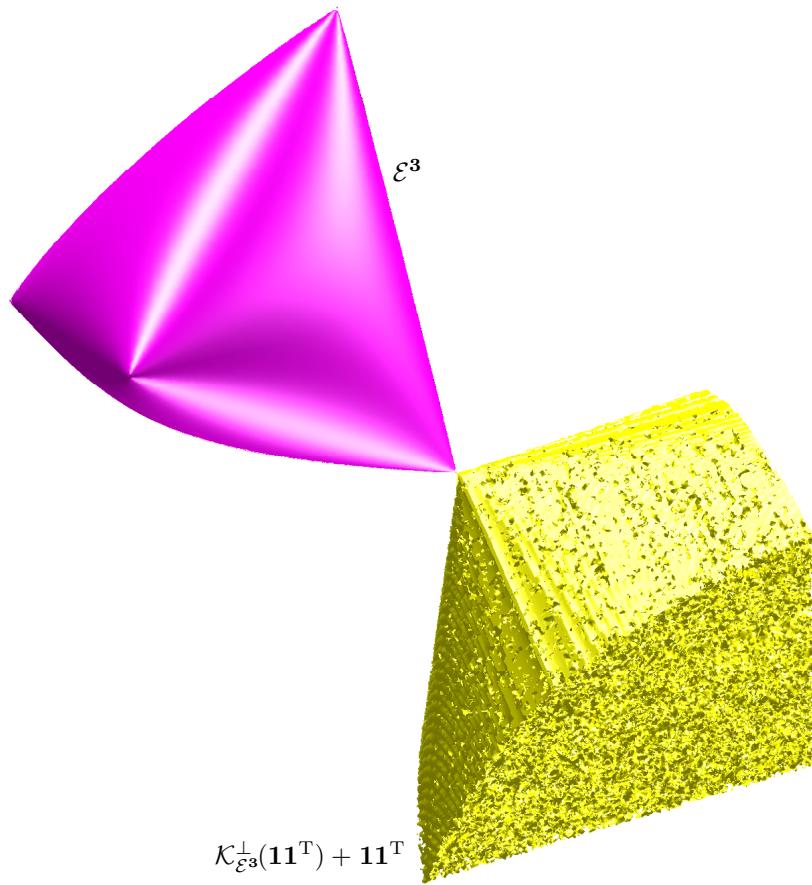
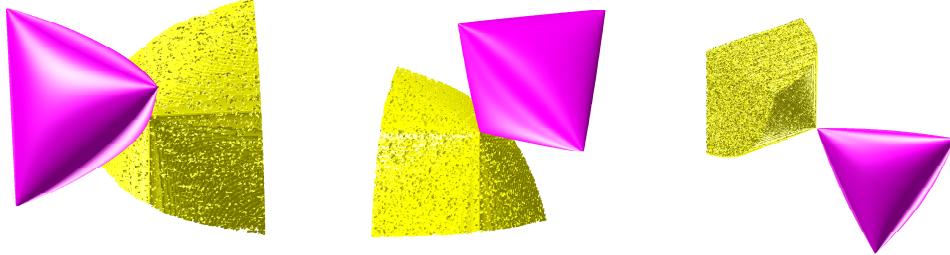


Figure 203: A few renderings (see next page) of normal cone  $\mathcal{K}_{\mathcal{E}^3}^\perp$  to ellipotope  $\mathcal{E}^3$  (Figure 155), at point  $\mathbf{1}\mathbf{1}^T$ , projected on  $\mathbb{R}^3$ . In [266, fig.2], normal cone is claimed circular in this dimension. (Numerical artifacts corrupt boundary and make truncated relative interior appear corporeal.)



**E.10.3.2.1 Definition.** *Normal cone.* [299] [43, p.261] [225, §A.5.2] (§2.13.11) [58, §2.1] [342, §3] [343, p.15] The normal cone to any set  $\mathcal{Z} \subseteq \mathbb{R}^n$  at any particular point  $a \in \mathbb{R}^n$  is defined as the closed cone

$$\begin{aligned}\mathcal{K}_{\mathcal{Z}}^\perp(a) &= \{z \in \mathbb{R}^n \mid z^T(y - a) \leq 0 \quad \forall y \in \mathcal{Z}\} = -(z - a)^* \\ &= \{z \in \mathbb{R}^n \mid z^T y \leq 0 \quad \forall y \in \mathcal{Z} - a\}\end{aligned}\quad (454)$$

an intersection of halfspaces about the origin in  $\mathbb{R}^n$ , hence convex regardless of convexity of  $\mathcal{Z}$ ; it is the negative dual cone to translate  $\mathcal{Z} - a$ ; the set of all vectors normal to  $\mathcal{Z}$  at  $a$  (§E.9.1.0.1).  $\triangle$

Examples of normal cone construction are illustrated in Figure 71, Figure 202, Figure 203, and Figure 204. Normal cone at  $\mathbf{0}$  in Figure 202 is the vector sum (§2.1.8) of two normal cones; [58, §3.3 exer.10] for  $\mathcal{H}_1 \cap \text{intr } \mathcal{H}_2 \neq \emptyset$

$$\mathcal{K}_{\mathcal{H}_1 \cap \mathcal{H}_2}^\perp(\mathbf{0}) = \mathcal{K}_{\mathcal{H}_1}^\perp(\mathbf{0}) + \mathcal{K}_{\mathcal{H}_2}^\perp(\mathbf{0}) \quad (2288)$$

This formula applies more generally to other points in the intersection.

The normal cone to any affine set  $\mathcal{A}$  at  $\alpha \in \mathcal{A}$ , for example, is the orthogonal complement of  $\mathcal{A} - \alpha$ . When  $\mathcal{A} = \mathbf{0}$ ,  $\mathcal{K}_{\mathcal{A}}^\perp(\mathbf{0}) = \mathcal{A}^\perp$  is  $\mathbb{R}^n$  the ambient space of  $\mathcal{A}$ .

Projection of any point in the translated normal cone  $\mathcal{K}_{\mathcal{C}}^\perp(a \in \mathcal{C}) + a$  on convex set  $\mathcal{C}$  is identical to  $a$ ; in other words, point  $a$  is that point in  $\mathcal{C}$  closest to any point belonging to the translated normal cone  $\mathcal{K}_{\mathcal{C}}^\perp(a) + a$ ; e.g., Theorem E.4.0.0.1. The normal cone to convex cone  $\mathcal{K}$  at the origin

$$\mathcal{K}_{\mathcal{K}}^\perp(\mathbf{0}) = -\mathcal{K}^* \quad (2289)$$

is the negative dual cone. Any point belonging to  $-\mathcal{K}^*$ , projected on  $\mathcal{K}$ , projects on the origin. More generally, [122, §4.5]

$$\mathcal{K}_{\mathcal{K}}^\perp(a) = -(K - a)^* \quad (2290)$$

$$\mathcal{K}_{\mathcal{K}}^\perp(a \in \mathcal{K}) = -\mathcal{K}^* \cap a^\perp \quad (2291)$$

The normal cone to  $\bigcap \mathcal{C}_k$  at  $Pb$  in Figure 196 is ray  $\{\xi(b - Pb) \mid \xi \geq 0\}$  illustrated in Figure 202. Applying Dykstra's algorithm to that example, convergence to the desired result is achieved in two iterations as illustrated in Figure 201. Yet applying

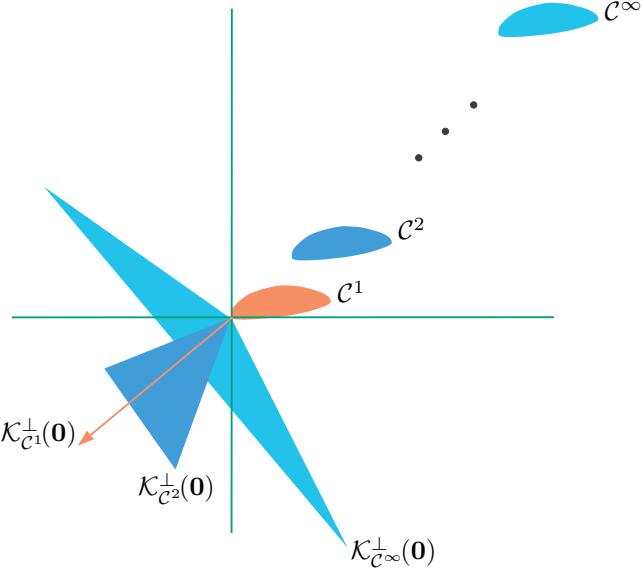


Figure 204: Rough sketch of normal cone to set  $\mathcal{C} \subset \mathbb{R}^2$  as  $\mathcal{C}$  wanders toward infinity. A point at which a normal cone is determined, here the origin, need not belong to the set. Normal cone to  $\mathcal{C}^1$  is a ray. But as  $\mathcal{C}$  moves outward, normal cone approaches a halfspace.

Dykstra's algorithm to the example in Figure 195 does not improve rate of convergence, unfortunately, because the given point  $b$  and all the alternating projections already belong to the translated normal cone at the vertex of intersection.

#### E.10.3.3 speculation

Dykstra's algorithm always converges at least as quickly as classical alternating projection, never slower [122], and it succeeds where alternating projection fails. Rate of convergence is wholly dependent on particular geometry of a given problem. From these few examples we surmise, unique minimum-distance projection on *blunt* (not sharp or acute, informally) full-dimensional polyhedral cones may be found by Dykstra's algorithm in few iterations. But total number of alternating projections, constituting those iterations, can never be less than number of convex sets.

## Appendix F

# Notation, Definitions, Glossary

$b$	scalar or column vector (italic $abcdefghijklm諾qrstuvwxyz$ )
$b_i$	$i^{\text{th}}$ entry of vector $b = [b_i, i=1 \dots n]$ or $i^{\text{th}} b$ vector from a list $\{b_j, j=1 \dots n\}$ or $i^{\text{th}}$ iterate of vector $b$
$b_{i:j}$	or $b(i:j)$ : truncated vector comprising $i^{\text{th}}$ through $j^{\text{th}}$ entry of vector $b$ (290)
$b_k(i:j)$	or $b_{i:j,k}$ : truncated vector comprising $i^{\text{th}}$ through $j^{\text{th}}$ entry of vector $b_k$
$b^T$	vector transpose or row vector
$b^H$	complex conjugate transpose $b^{*T}$
$A$	matrix (italic $ABCDEFGHIJKLMNPQRSTUVWXYZ$ )
$A^T$	Matrix transpose $[A_{ij}] \leftarrow [A_{ji}]$ is a linear operator. Regarding $A$ as a linear operator, $A^T$ is its adjoint.
$A^{-T}$	matrix transpose of inverse; and <i>vice versa</i> , $(A^{-1})^T = (A^T)^{-1}$ (confer p.487)
$A^{T_1}$	first of various transpositions of a cubix or quartix $A$ (p.549, p.553)
$A^{-1}$	inverse of matrix $A$
$A^\dagger$	Moore-Penrose pseudoinverse of matrix $A$ (§E)
$A_{ij}$	or $A(i,j)$ : $ij^{\text{th}}$ entry of matrix $A = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}^{1 \atop 2 \atop 3}$ or rank-one matrix $a_i a_j^T$ (§4.10)
$A(i,j)$	$A$ is a function of $i$ and $j$
$A_i$	$i^{\text{th}}$ matrix from a set or $i^{\text{th}}$ principal submatrix (1313) or $i^{\text{th}}$ iterate of $A$
$A(i,:)$	$i^{\text{th}}$ row of matrix $A$
$A(:,j)$	$j^{\text{th}}$ column of matrix $A$ [181, §1.1.8]
$A_{i:j,k:\ell}$	or $A(i:j, k:\ell)$ : submatrix; $i^{\text{th}}$ through $j^{\text{th}}$ row and $k^{\text{th}}$ through $\ell^{\text{th}}$ column of $A$

$\sqrt{\phantom{x}}$	positive square root
$\sqrt[x]{\phantom{x}}$	entrywise positive square root of vector $x$
$\sqrt[\ell]{\phantom{x}}$	positive $\ell^{\text{th}}$ root
$A^{1/2}$ and $\sqrt{A}$	$A^{1/2}$ is any matrix such that $A^{1/2}A^{1/2}=A$ . For $A \in \mathbb{S}_+^n$ , $\sqrt{A} \in \mathbb{S}_+^n$ is unique and $\sqrt{A}\sqrt{A}=A$ . [58, §1.2] (§A.5.1.4)
$\sqrt[3]{D}$	$= [\sqrt{d_{ij}}]$ absolute distance matrix (1501) or Hadamard positive square root: $D = \sqrt[3]{D} \circ \sqrt[3]{D}$
<i>thin</i>	a skinny matrix; meaning, more rows than columns: $\left[ \quad \right]$ . <i>When there are more equations than unknowns, we say that the system <math>Ax=b</math> is overdetermined.</i> [181, §5.3]
<i>wide</i>	a fat matrix; meaning, more columns than rows: $\left[ \quad \right]$ . <i>underdetermined</i>
<i>hollow</i>	matrix having <b>0</b> main diagonal
$\mathcal{A}$	some set (calligraphic $\mathcal{ABCDEFHIJKLMNOPQRSTUVWXYZ}$ )
$\mathbb{A}$	set of vectors or matrices (blackboard $\mathbb{ABCDEFHIJKLMNOPQRSTUVWXYZ}$ )
$\mathfrak{F}$	discrete Fourier transform (918) (Euler Fraktur $\mathfrak{A}\mathfrak{B}\mathfrak{C}\mathfrak{D}\mathfrak{E}\mathfrak{F}\mathfrak{G}\mathfrak{H}\mathfrak{I}\mathfrak{J}\mathfrak{K}\mathfrak{L}\mathfrak{M}\mathfrak{N}\mathfrak{O}\mathfrak{P}\mathfrak{Q}\mathfrak{R}\mathfrak{S}\mathfrak{T}\mathfrak{U}\mathfrak{V}\mathfrak{W}\mathfrak{X}\mathfrak{Y}\mathfrak{Z}$ )
$\mathcal{F}(\mathcal{C} \ni A)$	smallest face (174) that contains element $A$ of set $\mathcal{C}$
$\mathcal{G}(\mathcal{K})$	generators (§2.13.4.2.1) of set $\mathcal{K}$ ; any collection of points and directions whose hull constructs $\mathcal{K}$
$\mathcal{L}_\nu$	level set (562)
$\mathcal{L}_\nu$	sublevel set (566)
$\mathcal{L}^\nu$	superlevel set (655)
$\mathfrak{L}$	Lagrangian (512)
$\mathfrak{E}$	member of ellipope $\mathcal{E}_t$ (1242) parametrized by scalar $t$
$\mathcal{E}$	ellipope (1221)
$E$	elementary matrix
$E_{ij}$	member of standard orthonormal basis for symmetric (60) or symmetric hollow (76) matrices
<i>id est</i>	from the Latin meaning <i>that is</i>
<i>e.g.</i>	<i>exempli gratia</i> , from the Latin meaning <i>for sake of example</i>
<i>sic</i>	from the Latin meaning <i>so</i> or <i>thus</i> or <i>in this manner</i> ; something meant as written
<i>videlicet</i>	from the Latin meaning <i>it is permitted to see</i>
<i>ibidem</i>	from the Latin meaning <i>in the same place</i>

<i>no.</i>	<i>number</i> , from the Latin <i>numero</i>
a.i.	affinely independent ( <a href="#">§2.4.2.3</a> )
c.i.	conically independent ( <a href="#">§2.10</a> )
l.i.	linearly independent
w.r.t	<i>with respect to</i>
<b>a.k.a</b>	<i>also known as</i>
re	real part
im	imaginary part
$\imath$ or $\jmath$	$\sqrt{-1}$
$\subseteq$	<i>subset, superset</i>
$\subset$	<i>proper subset, proper superset</i>
$\cap$	<i>intersection, union</i>
$\in$	<i>membership, element belongs to, or element is a member of</i>
$\ni$	<i>membership, contains as in <math>\mathcal{C} \ni y</math> (<math>\mathcal{C}</math> contains element <math>y</math>)</i>
$\exists$	<i>such that</i>
$\exists$	<i>there exists</i>
$\therefore$	<i>therefore</i>
$\forall$	<i>for all, or over all</i>
&	(ampersand) <i>and</i>
$\&$	(ampersand italic) <i>and</i>
$\propto$	<i>proportional to</i>
$\infty$	<i>infinity</i>
$\equiv$	<i>equivalent to</i>
$\triangleq$	<i>defined equal to, equal by definition</i>
$\approx$	<i>approximately equal to</i>
$\simeq$	<i>isomorphic to or with</i>
$\cong$	<i>congruent to or with</i>
$\overline{\phantom{x}}$	Hadamard quotient as in, for $x, y \in \mathbb{R}^n$ , $\frac{x}{y} \triangleq [x_i/y_i, i=1 \dots n] \in \mathbb{R}^n$
$\circ$	Hadamard product of matrices: $x \circ y \triangleq [x_i y_i, i=1 \dots n] \in \mathbb{R}^n$ ( <a href="#">§D.1.2.2</a> , <a href="#">§A.1.1</a> )
$\otimes$	Kronecker product of matrices ( <a href="#">§D.1.2.1</a> , <a href="#">§A.1.1</a> )

- $\oplus$  vector sum of sets  $\mathcal{X} = \mathcal{Y} \oplus \mathcal{Z}$  where every element  $x \in \mathcal{X}$  has unique expression  $x = y + z$  where  $y \in \mathcal{Y}$  and  $z \in \mathcal{Z}$ ; [343, p.19] then summands are *algebraic complements*.  $\mathcal{X} = \mathcal{Y} \oplus \mathcal{Z} \Rightarrow \mathcal{X} = \mathcal{Y} + \mathcal{Z}$ . Now assume  $\mathcal{Y}$  and  $\mathcal{Z}$  are nontrivial subspaces.  $\mathcal{X} = \mathcal{Y} + \mathcal{Z} \Rightarrow \mathcal{X} = \mathcal{Y} \oplus \mathcal{Z} \Leftrightarrow \mathcal{Y} \cap \mathcal{Z} = \mathbf{0}$  [344, §1.2] [122, §5.8]. Each element from a vector sum (+) of subspaces has unique expression ( $\oplus$ ) when a basis from each subspace is linearly independent of bases from all the other subspaces.
- $\ominus$  likewise, unique vector difference of sets
- $\boxplus$  orthogonal vector sum of sets  $\mathcal{X} = \mathcal{Y} \boxplus \mathcal{Z}$  where every element  $x \in \mathcal{X}$  has unique orthogonal expression  $x = y + z$  where  $y \in \mathcal{Y}$ ,  $z \in \mathcal{Z}$ , and  $y \perp z$ . [366, p.51]  $\mathcal{X} = \mathcal{Y} \boxplus \mathcal{Z} \Rightarrow \mathcal{X} = \mathcal{Y} + \mathcal{Z}$ . If  $\mathcal{Z} \subseteq \mathcal{Y}^\perp$  then  $\mathcal{X} = \mathcal{Y} \boxplus \mathcal{Z} \Leftrightarrow \mathcal{X} = \mathcal{Y} \oplus \mathcal{Z}$ . [122, §5.8] If  $\mathcal{Z} = \mathcal{Y}^\perp$  then summands are *orthogonal complements*.
- $\pm$  plus or minus or plus and minus
- $\setminus$  as in  $\setminus \mathcal{A}$  means logical not  $\mathcal{A}$ , or relative complement of set  $\mathcal{A}$ ; id est,  $\setminus \mathcal{A} = \{x \notin \mathcal{A}\}$ ; e.g.,  $\mathcal{B} \setminus \mathcal{A} \triangleq \{x \in \mathcal{B} \mid x \notin \mathcal{A}\} \equiv \mathcal{B} \cap \setminus \mathcal{A}$
- $\Rightarrow$  and  $\Leftarrow$  sufficient and necessary, implies and is implied by; e.g,  
 $A$  is sufficient:  $A \Rightarrow B$ ,  $A$  is necessary:  $A \Leftarrow B$ ,  
 $A \Rightarrow B \Leftrightarrow \setminus A \Leftarrow \setminus B$ ,  $A \Leftarrow B \Leftrightarrow \setminus A \Rightarrow \setminus B$ ,  
if  $A$  then  $B$ , if  $B$  then  $A$ ,  
 $A$  only if  $B$ ,  $B$  only if  $A$ .
- $\Leftrightarrow$  if and only if (iff) or corresponds with or necessary and sufficient or logical equivalence
- $is$  as in  $A$  is  $B$  means  $A \Rightarrow B$ ; conventional usage of English language imposed by logicians
- $\not\Rightarrow$  and  $\not\Leftarrow$  insufficient and unnecessary, does not imply and is not implied by; e.g,  
 $A$  is insufficient:  $A \not\Rightarrow B$ ,  $A$  is unnecessary:  $A \not\Leftarrow B$ ,  
 $A \not\Rightarrow B \Leftrightarrow \setminus A \not\Leftarrow \setminus B$ ,  $A \not\Leftarrow B \Leftrightarrow \setminus A \not\Rightarrow \setminus B$ .
- $\leftarrow$  is replaced with; substitution, assignment
- $\rightarrow$  goes to, or approaches, or maps to
- $t \rightarrow 0^+$   $t$  goes to 0 from above; meaning, from the positive [225, p.2]
- $\vdots \quad \ddots \quad \dots$  as in  $1 \dots 1$  means ones in a row or  
 $[s_1 \dots s_N]$  means continuation; a matrix whose columns are  $s_i$  for  $i=1 \dots N$   
or as in  $n(n-1)(n-2) \dots 1$  means continuation of a product
- $\dots$  as in  $i=1 \dots N$  meaning,  $i$  is a sequence of successive integers beginning with 1 and ending with  $N$ ; id est,  $1 \dots N = 1:N$
- $:$  as in  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  meaning  $f$  is a mapping,  
or sequence of successive integers specified by bounds as in  $i:j = i \dots j$   
(if  $j < i$  then sequence is descending)
- $f$  real function or multidimensional function a.k.a operator
- $f : \mathcal{M} \rightarrow \mathcal{R}$  meaning  $f$  is a mapping from ambient space  $\mathcal{M}$  to ambient  $\mathcal{R}$ , not necessarily denoting either domain or range

	as in $f(x) \mid x \in \mathcal{C}$ means <i>with the condition(s)</i> or <i>such that</i> or <i>evaluated for</i> , or as in $\{f(x) \mid x \in \mathcal{C}\}$ means <i>evaluated for each and every</i> $x$ <i>belonging to set</i> $\mathcal{C}$
$g _{x_p}$	<i>expression g evaluated at</i> $x_p$
	<i>parallel</i>
$A, B$	as in, for example, $A, B \in \mathbb{S}^N$ means $A \in \mathbb{S}^N$ and $B \in \mathbb{S}^N$
$(a, b)$	<i>open interval</i> between $a$ and $b$ in $\mathbb{R}$ or <i>variable pair</i> perhaps of disparate dimension
$[a, b]$	<i>closed interval</i> or <i>line segment</i> between $a$ and $b$ in $\mathbb{R}$
( )	<i>hierachal, parenthetical, optional</i>
$\binom{n}{k}$	$\triangleq \begin{cases} 1, & k = 0 \\ -1 \frac{k(-n+k-1)!}{k!(-n-1)!}, & k > 0 \\ -1^{n-k} \frac{(-k-1)!}{(n-k)!(-n-1)!}, & k \leq n \\ 0, & n < k < 0 \\ 0, & k < 0 \text{ or } k > n \\ \frac{n!}{k!(n-k)!}, & \text{otherwise} \end{cases} \quad n < 0 \quad [255] \quad \text{binomial coefficient on } \mathbb{Z}^2$
!	<i>factorial; id est</i> , for integer $n > 0$ , $n! \triangleq n(n-1)(n-2)\cdots 1$ , $(-n)! \triangleq \infty$ , $0! \triangleq 1$
{ }	curly braces denote a <i>set</i> or <i>list</i> ; e.g., $\{Xa \mid a \succeq 0\}$ the <i>set comprising Xa evaluated for each and every</i> $a \succeq 0$ where membership of $a$ to some space is implicit, a <i>union</i> ; or $\{0, 1\}^n$ represents a binary vector of dimension $n$
$\langle \rangle$	angle brackets denote <i>vector inner-product</i> (33) (38)
[ ]	matrix or vector, or quote insertion, or citation
$[d_{ij}]$	matrix whose $ij^{\text{th}}$ entry is $d_{ij}$
$[x_i]$	vector whose $i^{\text{th}}$ entry is $x_i$
$x_p$	particular value of $x$
$x_0$	particular instance of $x$ , or initial value of a sequence $x_i$
$x_1$	first entry of vector $x$ , or first element of a list $\{x_i\}$
$x_\varepsilon$	<i>extreme point</i>
$x_+$	vector $x$ whose negative entries are replaced with 0: $x_+ = \frac{1}{2}(x +  x )$ <span style="float: right;">(540)</span> <i>nonnegative part of x or clipped vector x</i>
$x_-$	$x_- \triangleq \frac{1}{2}(x -  x )$ : <i>nonpositive part of x</i> $= x_+ + x_-$
$\check{x}$	known data
$x^*$	optimal value of variable $x$ . optimal $\Rightarrow$ feasible
$x^*$	<i>complex conjugate</i> or <i>dual variable</i> or <i>extreme direction of dual cone</i>
$f^*$	<i>convex conjugate function</i> $f^*(s) = \sup\{\langle s, x \rangle - f(x) \mid x \in \text{dom } f\}$
$P_C x$ or $Px$	projection of point $x$ on set $\mathcal{C}$ , $P$ is operator or idempotent matrix

$P_k x$	projection of point $x$ on set $\mathcal{C}_k$ or on range of implicit vector
$\delta(A)$	(a.k.a $\text{diag}(A)$ , §A.1) <i>vector made from main diagonal of <math>A</math></i> if $A$ is a matrix; otherwise, <i>diagonal matrix made from vector <math>A</math></i>
$\delta^2(A)$	$\equiv \delta(\delta(A))$ . For vector or diagonal matrix $\Lambda$ , $\delta^2(\Lambda) = \Lambda$
$\delta(A)^2$	$= \delta(A)\delta(A)$ where $A$ is a vector
$\lambda_i(X)$	$i^{\text{th}}$ entry of vector $\lambda$ is function of $X$
$\lambda(X)_i$	$i^{\text{th}}$ entry of vector-valued function of $X$
$\lambda(A)$	<i>vector of eigenvalues of matrix <math>A</math></i> , (1615) typically arranged in nonincreasing order
$\lambda(\mathbb{A})$	spectral cone $\mathcal{K}_\lambda$ for matrix set $\mathbb{A}$ (§5.11.2.3)
$\sigma(A)$	<i>vector of singular values of matrix <math>A</math></i> (always arranged in nonincreasing order), or <i>support function in direction <math>A</math></i>
$\Sigma$	diagonal matrix of singular values, not necessarily square
$\sum$	sum. Empty sum equals, conventionally, 0 or $\mathbf{0}$
$\pi(\gamma)$	nonlinear <i>permutation operator</i> (or <i>presorting function</i> ) arranges vector $\gamma$ into nonincreasing order (§7.1.3). $\pi_i$ is a permutation matrix; e.g., (947).
$\Xi$	permutation matrix
$\Pi$	doublet or permutation operator or matrix or set of all permutation matrices
$\prod$	product. Empty product equals, conventionally, 1 or $I$
$\psi(Z)$	signum-like <i>step function</i> that returns a scalar for matrix argument (754), it returns a vector for vector argument (1733)
$D$	symmetric hollow matrix of distance-square or <i>Euclidean distance matrix</i>
$\mathbf{D}$	Euclidean distance matrix operator
$\mathbf{D}^T(X)$	adjoint operator
$\mathbf{D}(X)^T$	transpose of $\mathbf{D}(X)$
$\mathbf{D}^{-1}(X)$	inverse operator
$\mathbf{D}(X)^{-1}$	inverse of $\mathbf{D}(X)$
$D^*$	optimal value of variable $D$
$D^*$	dual to variable $D$
$\mathbf{V}$	geometric centering operator, $\mathbf{V}(D) = -VDV^{\frac{1}{2}}$ (1141)
$\mathbf{V}_{\mathcal{N}}$	$\mathbf{V}_{\mathcal{N}}(D) = -V_{\mathcal{N}}^T D V_{\mathcal{N}}$ (1155)
$V$	$N \times N$ symmetric elementary, auxiliary, projector, geometric centering matrix, $\mathcal{R}(V) = \mathcal{N}(\mathbf{1}^T)$ , $\mathcal{N}(V) = \mathcal{R}(\mathbf{1})$ , $V^2 = V$ (§B.4.1)
$V_{\mathcal{N}}$	$N \times N-1$ Schoenberg auxiliary matrix $\mathcal{R}(V_{\mathcal{N}}) = \mathcal{N}(\mathbf{1}^T)$ , $\mathcal{N}(V_{\mathcal{N}}^T) = \mathcal{R}(\mathbf{1})$ (§B.4.2)

$V_{\mathcal{X}}$	$V_{\mathcal{X}}V_{\mathcal{X}}^T \equiv V^T X^T X V$ (1333)
$X$	point list ((77) having cardinality $N$ ) arranged columnar in $\mathbb{R}^{n \times N}$ , or set of generators, or extreme directions, or matrix variable
$G$	Gram matrix $X^T X$ (1042)
$r$	affine dimension (1187)
$\alpha_c$	geometric center (1118)
$\rho$	rank of matrix or bound on affine dimension
$k$	number of conically independent generators
$\mathbb{k}$	raw-data domain of Magnetic Resonance Imaging (MRI) machine, as in $\mathbb{k}$ -space
$n$	Euclidean (ambient spatial) dimension of list $X \in \mathbb{R}^{n \times N}$ , or integer
$\eta$	noise factor or noise signal or normal vector or $\min\{m, n\}$
$N$	cardinality of list $X \in \mathbb{R}^{n \times N}$ , or integer
epi	function epigraph
dom	function domain
$\mathcal{R}f$	function range
$\mathcal{R}(A^T)$	the subspace: <i>rowspace of A</i> (143) or span basis $\mathcal{R}(A^T)$ ; $\mathcal{R}(A^T) \perp \mathcal{N}(A)$
$\mathcal{R}(A)$	the subspace: <i>range of A</i> (144) or span basis $\mathcal{R}(A)$ ; $\mathcal{R}(A) \perp \mathcal{N}(A^T)$
span	as in $\text{span } A = \mathcal{R}(A) = \{Ax \mid x \in \mathbb{R}^n\}$ when $A$ is a matrix
basis $\mathcal{R}(A)$	<i>overcomplete columnar basis for range of A</i> or <i>minimal set constituting generators for vertex-description of <math>\mathcal{R}(A)</math></i> or <i>linearly independent set of vectors spanning <math>\mathcal{R}(A)</math></i>
$\mathcal{N}(A)$	the subspace: <i>nullspace of A</i> (145) a.k.a kernel of $A$ ; $\mathcal{N}(A) \perp \mathcal{R}(A^T)$
$\mathbb{R}^n$	Euclidean $n$ -dimensional real vector space (nonnegative integer $n$ ); a subspace, conventionally, but not a proper subspace. [254, §2.1] $\mathbb{R}^0 = \mathbf{0}$ . $\mathbb{R} = \mathbb{R}^1$ or vector space of unspecified dimension. [442]
$\mathbb{R}^{m \times n}$	Euclidean vector space of $m$ by $n$ dimensional real matrices
$\times$	Cartesian product. $\mathbb{R}^{m \times n-m} \triangleq \mathbb{R}^{m \times (n-m)}$ . $\mathcal{K}_1 \times \mathcal{K}_2 = \begin{bmatrix} \mathcal{K}_1 \\ \mathcal{K}_2 \end{bmatrix}$
$\begin{bmatrix} \mathbb{R}^m \\ \mathbb{R}^n \end{bmatrix}$	$\mathbb{R}^m \times \mathbb{R}^n = \mathbb{R}^{m+n}$
$\mathbb{Z}$	the real integers
$\mathbb{N}$	the nonnegative natural numbers; <i>id est</i> , $\mathbb{Z}_+$
$\mathbb{B}^n$ , $\mathbb{B}^{n \times n}$	$\{0, 1\}^n$ and $\{0, 1\}^{n \times n}$ binary vectors of respective dimension $n$ and $n \times n$
$\mathbb{B}_{\pm}^n$ , $\mathbb{B}_{\pm}^{n \times n}$	$\{-1, 1\}^n$ and $\{-1, 1\}^{n \times n}$ bipolar binary vectors of dimension $n$ and $n \times n$
$\mathbb{C}^n$ , $\mathbb{C}^{n \times n}$	Euclidean complex vector space of respective dimension $n$ and $n \times n$

$\mathbb{R}_+^n$ , $\mathbb{R}_+^{n \times n}$	nonnegative orthant in Euclidean vector space of respective dimension $n$ and $n \times n$
$\mathbb{R}_-^n$ , $\mathbb{R}_-^{n \times n}$	nonpositive orthant in Euclidean vector space of respective dimension $n$ and $n \times n$
$\mathbb{S}^n$	subspace of real symmetric $n \times n$ matrices; the <i>symmetric matrix subspace</i> . $\mathbb{S}^0 = \mathbf{0}$ . $\mathbb{S} = \mathbb{S}^1$ or symmetric subspace of unspecified dimension.
$\mathbb{S}^{n \perp}$	orthogonal complement of $\mathbb{S}^n$ in $\mathbb{R}^{n \times n}$ , the antisymmetric matrices (52)
$\mathbb{S}_+^n$	convex cone comprising all (real) symmetric positive semidefinite $n \times n$ matrices, the <i>positive semidefinite cone</i>
$\text{intr } \mathbb{S}_+^n$	interior of convex cone comprising all (real) symmetric positive semidefinite $n \times n$ matrices; <i>id est</i> , positive definite matrices
$\mathbb{S}_+^n(\rho)$	$= \{X \in \mathbb{S}_+^n \mid \text{rank } X \geq \rho\}$ (265) convex set of all positive semidefinite $n \times n$ symmetric matrices whose rank equals or exceeds $\rho$
$\mathbf{EDM}^N$	cone of $N \times N$ Euclidean distance matrices in the symmetric hollow subspace
$\sqrt{\mathbf{EDM}^N}$	nonconvex cone of $N \times N$ Euclidean absolute distance matrices in the symmetric hollow subspace (§6.3)
$\mathbb{S}_0^n$	subspace comprising all symmetric $n \times n$ matrices having all zeros in first row and column (2200) (§5.4.2.1)
$\mathbb{S}_h^n$	subspace comprising all symmetric hollow $n \times n$ matrices ( $\mathbf{0}$ main diagonal), the <i>symmetric hollow subspace</i> (67)
$\mathbb{S}_h^{n \perp}$	orthogonal complement of $\mathbb{S}_h^n$ in $\mathbb{S}^n$ , the set of all diagonal matrices (68)
$\mathbb{S}_c^n$	subspace comprising all geometrically centered symmetric $n \times n$ matrices; <i>geometric center subspace</i> $\mathbb{S}_c^N = \{Y \in \mathbb{S}^N \mid Y\mathbf{1} = \mathbf{0}\}$ (2196)
$\mathbb{S}_c^{n \perp}$	orthogonal complement of $\mathbb{S}_c^n$ in $\mathbb{S}^n$ (2198); <i>translation-invariant subspace</i>
$\mathbb{R}_c^{m \times n}$	subspace comprising all geometrically centered $m \times n$ matrices (2195)
$\mathbb{R}_h^{n \times n}$	subspace of symmetric [sic] matrices having $\mathbf{0}$ main diagonal; a.k.a., <i>real hollow subspace</i> (64)
$\mathbb{R}_h^{n \times n \perp}$	subspace of antisymmetric antihollow matrices (65)
$X^\perp$	basis $\mathcal{N}(X^T)$ (§2.13.10, §E.3.4)
$x^\perp$	$\mathcal{N}(x^T); \{y \in \mathbb{R}^n \mid x^T y = 0\}$ (§2.13.11.1.1)
$\perp$	as in $A \perp B$ meaning $A$ is <i>orthogonal to</i> $B$ (and <i>vice versa</i> ), where $A$ and $B$ are sets, vectors, or matrices. When $A$ and $B$ are vectors (33) (34) (or matrices (38) under Frobenius' norm), $A \perp B \Leftrightarrow \langle A, B \rangle = 0 \Leftrightarrow \ A + B\ ^2 = \ A\ ^2 + \ B\ ^2$
$\mathcal{R}(P)^\perp$	$\mathcal{N}(P^T)$ ; orthogonal complement of $\mathcal{R}(P)$ (fundamental subspace relations (139))
$\mathcal{N}(P)^\perp$	$\mathcal{R}(P^T)$
$\mathcal{R}^\perp$	$= \{y \in \mathbb{R}^n \mid \langle x, y \rangle = 0 \forall x \in \mathcal{R}\}$ (377). <i>Orthogonal complement of <math>\mathcal{R}</math> in <math>\mathbb{R}^n</math> when <math>\mathcal{R}</math> is a subspace</i>
$\mathcal{K}^\perp$	normal cone (454)

$\mathcal{A}^\perp$	normal cone to affine subset $\mathcal{A}$	(§3.1.1.2.2)
$\mathcal{K}$	cone	
$\mathcal{K}^*$	dual cone $-\mathcal{K}^\circ$	
$\mathcal{K}^\circ$	polar cone $-\mathcal{K}^*$	
$D^\circ$	polar variable $D$	
$360^\circ$	angular <i>degree</i> ; <i>e.g.</i> , $360^\circ \Leftrightarrow 2\pi$ radians	
$\mathcal{K}_{\mathcal{M}+}$	monotone nonnegative cone	
$\mathcal{K}_{\mathcal{M}}$	monotone cone	
$\mathcal{K}_\lambda$	spectral cone	
$\mathcal{K}_{\lambda\delta}^*$	cone of majorization	
$\mathcal{H}$	halfspace	
$\mathcal{H}_-$	halfspace described using an outward-normal (108) to the hyperplane partially bounding it	
$\mathcal{H}_+$	halfspace described using an inward-normal (109) to the hyperplane partially bounding it	
$\partial\mathcal{H}$	hyperplane; <i>id est</i> , partial boundary of halfspace	
$\underline{\partial\mathcal{H}}$	supporting hyperplane	
$\overline{\partial\mathcal{H}}_-$	a supporting hyperplane having outward-normal with respect to set it supports	
$\overline{\partial\mathcal{H}}_+$	a supporting hyperplane having inward-normal with respect to set it supports	
$\partial$	<i>partial derivative</i> or <i>partial differential</i> or <i>matrix of distance-square squared</i> (1541) or <i>boundary</i> of set $\mathcal{K}$ as in $\partial\mathcal{K}$ (17) (24)	
$d$	<i>derivative</i> or <i>differential</i>	
$\sqrt{d_{ij}}$	(absolute) distance scalar	
$d_{ij}$	distance-square scalar, EDM entry	
$\underline{d}$	vector of distance-square	
$\underline{d}_{ij}$	lower bound on distance-square $d_{ij}$	
$\overline{d}_{ij}$	upper bound on distance-square $d_{ij}$	
$\overline{AB}$	closed line segment between points A and B	
$AB$	matrix multiplication of $A$ and $B$	
$\bar{\mathcal{C}}$	<i>closure of set</i> $\mathcal{C}$	
$g'$	first derivative of possibly multidimensional function with respect to real argument	
$g''$	second derivative with respect to real argument	

$\overset{\rightarrow}{dg}^Y$	first directional derivative of possibly multidimensional function $g$ in direction $Y \in \mathbb{R}^{K \times L}$ (maintains dimensions of $g$ )
$\overset{\rightarrow}{dg}^2$	second directional derivative of $g$ in direction $Y$
$\nabla$	gradient from calculus, $\nabla f$ is shorthand for $\nabla_x f(x)$ . $\nabla f(y)$ means $\nabla_y f(y)$ or gradient $\nabla_x f(y)$ of $f(x)$ with respect to $x$ evaluated at $y$
$\nabla^2$	second-order gradient
$\Delta$	difference or discrete differential or distance scalar (Figure 28) or first-order difference matrix (931) or infinitesimal difference operator (§D.1.4)
$\triangle_{ijk}$	triangle made by vertices $i$ , $j$ , and $k$
$\ddot{v}$	coefficient vector for two spectral factors, Figure 110 level 2
$\ddot{\ddot{v}}$	coefficient vector corresponding to four spectral factors, Figure 110 level 3
$\ddot{\ddot{\ddot{v}}}$	vector containing numerator or denominator coefficients of eight spectral factors; level 4 in a bifurcation tree like Figure 110
CPU	central processing unit
dB	decibel
DC	direct current (0 Hz)
DCT	discrete cosine transform
DFT	discrete Fourier transform $\mathfrak{F}$ (918)
Hz	hertz (cycles per second), kHz means kilohertz, MHz megahertz, GHz gigahertz
EDM	Euclidean distance matrix
PSD	positive semidefinite
SDP	semidefinite program
LP	linear program
QUBO	quadratic unconstrained binary optimization
SVD	singular value decomposition
SNR	signal to noise ratio
USA	United States of America
<i>in</i>	function $f$ in $x$ means $x$ as argument to $f$ or $x$ in $\mathcal{C}$ means element $x$ is a member of set $\mathcal{C}$
<i>on</i>	function $f(x)$ on $\mathcal{A}$ means $\mathcal{A}$ is $\text{dom } f$ or relation $\preceq$ on $\mathcal{A}$ means $\mathcal{A}$ is set whose elements are subject to $\preceq$ or projection of $x$ on $\mathcal{A}$ means $\mathcal{A}$ is body on which projection is made or operating on vector identifies argument type to $f$ as “vector”
<i>over</i>	function $f(x)$ over $\mathcal{C}$ means $f$ evaluated at each and every element of set $\mathcal{C}$

<i>one-to-one</i>	injective map or unique correspondence between sets
<i>onto</i>	<i>function <math>f(x)</math> maps onto <math>\mathcal{M}</math></i> means $f$ over its domain is a surjection w.r.t $\mathcal{M}$
<i>injection</i>	$f$ that is <i>one-to-one</i>
<i>surjection</i>	$f$ that is <i>onto</i>
<i>bijection</i>	$f$ that is <i>one-to-one</i> and <i>onto</i>
<i>orthant</i>	generalization of two-dimensional <i>quadrant</i> $\sqsubset$ to higher dimension
<i>orthogonality</i>	generalization of two-dimensional <i>perpendicularity</i> $\perp$ to higher dimension
<i>decomposition</i>	<i>orthonormal</i> (2097, p.576), <i>biorthogonal</i> (2072, p.571)
<i>expansion</i>	<i>orthogonal</i> (2107, p.577), <i>biorthogonal</i> (409, p.149)
<i>vector</i>	<i>column vector</i> in $\mathbb{R}^n$ ; identifiable by Cartesian coordinates of point at its head
<i>entry</i>	<i>scalar element or real variable constituting a vector or matrix</i>
<i>cubix</i>	member of $\mathbb{R}^{M \times N \times L}$
<i>quartix</i>	member of $\mathbb{R}^{M \times N \times L \times K}$
<i>feasible</i>	as in <i>feasible solution</i> , means satisfies the (“subject to”) constraints of an optimization problem, may or may not be optimal
<i>feasible set</i>	most simply, <i>the set of all variable values satisfying all constraints of an optimization problem</i>
<i>active set</i>	an inequality constraint is termed <i>active</i> when it is met with equality; the set of all active constraints
<i>solution set</i>	most simply, <i>the set of all optimal solutions to an optimization problem</i> ; a subset of the feasible set and not necessarily a single point
<i>set</i>	collection of elements in which order and multiplicity are ignored. A set member is called <i>element</i> [436]
<i>list</i>	ordered set retaining multiplicity
<i>optimal</i>	as in <i>optimal solution</i> , means a solution to an optimization problem. An optimal solution is not necessarily unique, but there is no better solution. $\text{optimal} \Rightarrow \text{feasible}$
<i>optimum</i>	optimal value, usually the objective. Can be unique
<i>same</i>	as in <i>same problem</i> , means optimal solution set for one problem is identical to optimal solution set of another (without transformation)
<i>equivalent</i>	as in <i>equivalent problem</i> , means optimal solution to one problem can be derived from optimal solution to another via suitable transformation
<i>convex</i>	as in <i>convex problem</i> , essentially means a convex objective function optimized over a convex set (§4)
<i>objective</i>	the three objectives of Optimization are <i>minimize</i> (not min), <i>maximize</i> (not max), and <i>find</i>

<i>program</i>	<i>Semidefinite program</i> is any convex minimization, maximization, or feasibility problem constraining a variable to a subset of a positive semidefinite cone. <i>Prototypical semidefinite program</i> conventionally means: a semidefinite program having linear objective, affine equality constraints, but no inequality constraints except for cone membership. (§4.1.1) <i>Linear program</i> is any feasibility problem, or minimization or maximization of a linear objective, constraining the variable to some polyhedron. (§2.13.1.1.1) <i>Prototypical linear program</i> conventionally means: a linear program having linear objective, affine equality constraints, but no inequality constraints except for membership to a nonnegative orthant. (§4.1)
<i>natural order</i>	with reference to stacking columns in a vectorization means <i>a vector made from superposing column 1 on top of column 2 then superposing the result on column 3 and so on; as in a vector made from entries of the main diagonal</i> $\delta(A)$ means <i>taken from left to right and top to bottom</i>
<i>partial order</i>	relation $\preceq$ is a partial order, on a set, if it possesses reflexivity, antisymmetry, and transitivity (§2.7.2.2)
<i>operator</i>	mapping to a vector space (a multidimensional function)
<i>projector</i>	short for <i>projection operator</i> ; not necessarily minimum-distance or represented by a matrix
<i>sparsity</i>	ratio of number of nonzero entries to matrix-dimension product
<i>tight</i>	with reference to a bound means <i>a bound that can be met</i> , with reference to an inequality means <i>equality is achievable</i>
<i>trivial</i>	with reference to <b>0</b> matrix, function, solution, or $\{0\}$ subspace
$\emptyset$	<i>empty set</i> , an implicit member of every set
0	real zero
<b>0</b>	<i>origin</i> or vector or matrix of zeros
<i>O</i>	<i>sort-index matrix</i>
<i>O</i>	<i>order of magnitude</i> or <i>polynomial order</i> or <i>computational intensity</i> : $O(N)$ is first-order, $O(N^2)$ is second-, and so on
1	real one
<b>1</b>	vector of ones. $\mathbf{1} = \delta^2(\mathbf{1})$ , $\delta(\mathbf{1}) = I$
$\mathbf{1}_m$	$\mathbf{1} \in \mathbb{R}^m$
$e_i$	vector whose $i^{\text{th}}$ entry is 1 (otherwise 0); $i^{\text{th}}$ member of the standard basis for $\mathbb{R}^m$ (61)
I	Roman numeral one or capital i
<i>I</i>	Identity operator or matrix $I = \delta^2(I)$ , $\delta(I) = \mathbf{1}$
$I_m$	$I \in \mathbb{S}^m$
$\mathcal{I}$	<i>index set</i> , a discrete set of indices

$\max$	<i>maximum</i> [225, §0.1.1] or <i>largest element of a totally ordered set</i>
<i>maximal</i>	characterizes a maximum that is, somehow, <i>not necessarily</i> unique; <i>id est</i> , maximum $\not\Rightarrow$ unique maximum
$\underset{x}{\text{maximize}}$	<i>find maximum of objective function w.r.t independent variables <math>x</math>.</i> Subscript $x \leftarrow x \in \mathcal{C}$ may hold implicit constraints if context clear; <i>e.g.</i> , semidefiniteness
$\arg$	<i>argument of operator or function, or variable of optimization</i>
$\sup \mathcal{X}$	<i>supremum of totally ordered set <math>\mathcal{X}</math>, least upper bound</i> , may or may not belong to set [225, §0.1.1]; <i>e.g.</i> , range $\mathcal{X}$ of real function
$\arg \sup_x f(x)$	<i>argument <math>x</math> at supremum of function <math>f</math></i> ; not necessarily unique or a member of function domain
<i>subject to</i>	specifies constraints of an optimization problem; generally, inequalities and affine equalities. <i>Subject to</i> implies: anything not an independent variable is constant, an assignment, or substitution
$\min$	<i>minimum</i> [225, §0.1.1] or <i>smallest element of a totally ordered set</i>
<i>minimal</i>	describes a minimum that is, in some sense, <i>not necessarily</i> unique; <i>id est</i> , minimum $\not\Rightarrow$ unique minimum
$\underset{x}{\text{minimize}}$	<i>find objective function minimum w.r.t independent variables <math>x</math>.</i> Subscript $x \leftarrow x \in \mathcal{C}$ may hold implicit constraints if context clear; <i>e.g.</i> , semidefiniteness
$\underset{x}{\text{find}}$	<i>find any feasible solution, specified by the (“subject to”) constraints, w.r.t independent variables <math>x</math>.</i> Subscript $x \leftarrow x \in \mathcal{C}$ may hold implicit constraints if context clear; <i>e.g.</i> , semidefiniteness. “find” denotes a <i>feasibility problem</i> ; it is the third objective of Optimization
$\inf \mathcal{X}$	<i>infimum of totally ordered set <math>\mathcal{X}</math>, greatest lower bound</i> , may or may not belong to set [225, §0.1.1]; <i>e.g.</i> , range $\mathcal{X}$ of real function
$\arg \inf_x f(x)$	<i>argument <math>x</math> at infimum of function <math>f</math></i> ; not necessarily unique or a member of function domain
<i>iff</i>	<i>if and only if, necessary and sufficient</i> ; also the meaning indiscriminately attached to appearance of the word “if” in the statement of a mathematical definition, [146, p.106] [289, p.4] an esoteric practice worthy of abolition because of ambiguity thus conferred
<i>rel</i>	relative
<i>intr</i>	interior
<i>lim</i>	limit
<i>sgn</i>	signum function or <i>sign</i> ; for $x \in \mathbb{R}^n$ , $\text{sgn}(x) = \begin{cases} x_i /  x_i , & x_i \neq 0 \\ 0, & x_i = 0 \end{cases}$
<i>round</i>	round to nearest integer
<i>mod</i>	modulus function

tr	matrix trace
rank	as in $\text{rank } A$ , <i>rank of matrix <math>A</math></i> ; $\dim \mathcal{R}(A)$
dim	dimension, $\dim \mathbb{R}^n = n$ , $\dim \mathbb{R}^{m \times n} = m \times n$ $\dim(x \in \mathbb{R}^n) = n$ , $\dim \mathcal{R}(x \in \mathbb{R}^n) = 1$ $\dim(A \in \mathbb{R}^{m \times n}) = m \times n$ , $\dim \mathcal{R}(A \in \mathbb{R}^{m \times n}) = \text{rank } A$
aff	affine hull
dim aff	affine dimension $r$
card	cardinality, <i>number of nonzero entries</i> $\text{card } x \triangleq \ x\ _0$ or $N$ is cardinality of list $X \in \mathbb{R}^{n \times N}$ (p.259)
conv	convex hull (§2.3.2)
cone	conic hull (§2.3.3)
cenv	convex envelope (§7.2.2.1)
content	of high-dimensional bounded polyhedron, volume in $\mathbb{R}^3$ , area in $\mathbb{R}^2$ , and so on
cof	matrix of cofactors corresponding to matrix argument
dist	absolute distance between point or set arguments; <i>e.g.</i> , $\text{dist}(x, \mathcal{B})$
vec	columnar vectorization of $m \times n$ matrix, Euclidean dimension $mn$ (37)
svec	columnar vectorization of symmetric $n \times n$ matrix, Euclidean dimension $n(n+1)/2$ (57)
dvec	columnar vectorization of symmetric hollow $n \times n$ matrix, Euclidean dimension $n(n-1)/2$ (74)
$\measuredangle(x, y)$	<i>complex sinusoid phase</i> or <i>angle</i> between vectors $x$ and $y$ , or <i>dihedral angle</i> between affine subsets
$\succeq$	generalized inequality; <i>e.g.</i> , $A \succeq 0$ means:
	<ul style="list-style-type: none"> <li>• <i>vector or matrix <math>A</math> must be expressible in a biorthogonal expansion having nonnegative coordinates with respect to extreme directions of some implicit pointed closed convex cone <math>\mathcal{K}</math></i> (§2.13.2.0.1, §2.13.8.1.1),</li> <li>• <b>or</b> <i>comparison to the origin with respect to some implicit pointed closed convex cone</i> (2.7.2.2),</li> <li>• <b>or</b> (when <math>\mathcal{K} = \mathbb{S}_+^n</math>) <i>matrix <math>A</math> belongs to the positive semidefinite cone of symmetric matrices</i> (nonnegative eigenvalues, §2.9.0.1),</li> <li>• <b>or</b> (when <math>\mathcal{K} = \mathbb{R}_+^n</math>) <i>vector <math>A</math> belongs to the nonnegative orthant</i> (each vector entry is nonnegative, §2.3.1.1)</li> </ul>
$\succeq_{\mathcal{K}}$	as in $x \succeq_{\mathcal{K}} z$ means $x - z \in \mathcal{K}$ (185)
$\succ$	strict generalized inequality, membership to cone interior; $A \succ 0$ means:
	<ul style="list-style-type: none"> <li>• <i>vector or matrix <math>A</math> must be expressible in a biorthogonal expansion having positive coordinates with respect to extreme directions of some implicit pointed closed convex cone <math>\mathcal{K}</math></i> (§2.13.2.0.1, §2.13.8.1.1),</li> </ul>

- or comparison to the origin with respect to the interior of some implicit pointed closed convex cone (2.7.2.2),
- or (when  $\mathcal{K} = \mathbb{S}_+^n$ ) matrix  $A$  belongs to the interior of the positive semidefinite cone of symmetric matrices (positive eigenvalues, §2.9.0.1),
- or (when  $\mathcal{K} = \mathbb{R}_+^n$ ) vector  $A$  belongs to the interior of the nonnegative orthant (each vector entry is positive, §2.3.1.1)

$\neq$	not positive definite
$\geq$	scalar inequality, greater than or equal to; comparison of scalars, or entrywise comparison of vectors or matrices with respect to $\mathbb{R}_+$
nonnegative	for $a \in \mathbb{R}^n$ , $a \succeq 0$ ; id est, nonnegative entries when w.r.t nonnegative orthant; coefficients of vector on boundary of or interior to pointed closed convex cone $\mathcal{K}$
$>$	greater than
positive	for $a \in \mathbb{R}^n$ , $a \succ 0$ ; id est, positive (nonzero) entries when w.r.t nonnegative orthant; coefficients to no vector on boundary of pointed closed convex cone $\mathcal{K}$
$\lfloor \cdot \rfloor$	floor function, $\lfloor x \rfloor$ is greatest integer not exceeding $x$
$\text{\small \fbox{}}$	download button
$  \cdot  $	entrywise absolute value of scalars, vectors, and matrices
log	natural (or Napierian) logarithm
det	matrix determinant
$\ x\ $	$= \sqrt{\sum_{j=1}^n  x_j ^2}$ Euclidean norm or vector 2-norm $\ x\ _2$ (§3.2)
$\ x\ _2^2$	$= x^T x = \langle x, x \rangle$ Euclidean norm square (§3.1.1.1)
$\ x\ _\ell$	$= \sqrt[\ell]{\sum_{j=1}^n  x_j ^\ell}$ vector $\ell$ -norm for $\ell \geq 1$ (convex) $\triangleq \sum_{j=1}^n  x_j ^\ell$ vector $\ell$ -norm for $0 \leq \ell < 1$ (violates §3.2 no.3)
$\ x\ _0$	$= \mathbf{1}^T  x ^0$ ( $0^0 \triangleq 0$ ) “0-norm” or cardinality of vector $x$ (card $x$ ) (§4.6.1)
$\ x\ _1$	$= \mathbf{1}^T  x $ 1-norm, dual infinity-norm (§3.2)
$\ x\ _\infty$	$= \max\{ x_j  \forall j\}$ infinity-norm (§3.2)
$\ x\ _k$	$= \sum_{i=1}^k \pi( x )_i$ $k$ -largest norm (§3.2.2.1)
$\ X\ _2$	$= \sup_{\ a\ =1} \ Xa\ _2 = \sigma_1 = \sqrt{\lambda(X^T X)_1}$ matrix 2-norm or spectral norm, (596) largest singular value [181, p.56]. For $x$ a vector: $\ \delta(x)\ _2 = \ x\ _\infty$ . $\ Xa\ _2 \leq \ X\ _2 \ a\ _2 \quad (2292)$
$\ X\ _2^*$	$= \mathbf{1}^T \sigma(X)$ nuclear norm, dual spectral norm (§C.2)
$\ X\ $	$= \sqrt{\sum_{i,j} X_{ij}^2}$ Frobenius' matrix norm $\ X\ _F$ (§2.2.1)



# Bibliography

- [1] Edwin A. Abbott. *Flatland: A Romance of Many Dimensions*. Seely & Co., London, sixth edition, 1884.
- [2] Bob Adams. High Performance 16-/18-Bit  $\Sigma\Delta$  Stereo ADCs: AD1878/AD1879. *Analog Devices Corp*, 1993.  
<http://www.convexoptimization.com/TOOLS/AD1879.pdf>
- [3] Suliman Al-Homidan and Henry Wolkowicz. Approximate and exact completion problems for Euclidean distance matrices using semidefinite programming. *Linear Algebra and its Applications*, 406:109–141, September 2005.  
<http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.165.5577>
- [4] Faiz A. Al-Khayyal and James E. Falk. Jointly constrained biconvex programming. *Mathematics of Operations Research*, 8(2):273–286, May 1983.  
<http://www.convexoptimization.com/TOOLS/Falk.pdf>
- [5] Abdo Y. Alfakih. On the uniqueness of Euclidean distance matrix completions. *Linear Algebra and its Applications*, 370:1–14, 2003.
- [6] Abdo Y. Alfakih. On the uniqueness of Euclidean distance matrix completions: the case of points in general position. *Linear Algebra and its Applications*, 397:265–277, 2005.
- [7] Abdo Y. Alfakih, Amir Khandani, and Henry Wolkowicz. Solving Euclidean distance matrix completion problems via semidefinite programming. *Computational Optimization and Applications*, 12(1):13–30, January 1999.  
<http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.14.8307>
- [8] Abdo Y. Alfakih and Henry Wolkowicz. On the embeddability of weighted graphs in Euclidean spaces. Research Report CORR 98-12, Department of Combinatorics and Optimization, University of Waterloo, May 1998.  
<http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.50.1160>  
Erratum: p.376 herein.
- [9] Abdo Y. Alfakih and Henry Wolkowicz. Matrix completion problems. In Henry Wolkowicz, Romesh Saigal, and Lieven Vandenberghe, editors, *Handbook of Semidefinite Programming: Theory, Algorithms, and Applications*, chapter 18. Kluwer, 2000.  
<http://www.convexoptimization.com/TOOLS/Handbook.pdf>
- [10] Abdo Y. Alfakih and Henry Wolkowicz. Two theorems on Euclidean distance matrices and Gale transform. *Linear Algebra and its Applications*, 340:149–154, 2002.
- [11] Farid Alizadeh. *Combinatorial Optimization with Interior Point Methods and Semi-Definite Matrices*. PhD thesis, University of Minnesota, Computer Science Department, Minneapolis Minnesota USA, October 1991.
- [12] Farid Alizadeh. Interior point methods in semidefinite programming with applications to combinatorial optimization. *SIAM Journal on Optimization*, 5(1):13–51, February 1995.
- [13] Kurt Anstreicher and Henry Wolkowicz. On Lagrangian relaxation of quadratic matrix constraints. *SIAM Journal on Matrix Analysis and Applications*, 22(1):41–55, 2000.
- [14] Howard Anton. *Elementary Linear Algebra*. Wiley, second edition, 1977.
- [15] James Aspnes, David Goldenberg, and Yang Richard Yang. On the computational complexity of sensor network localization. In *Proceedings of the First International Workshop on Algorithmic Aspects of Wireless Sensor Networks (ALGOSENSORS)*, volume 3121 of *Lecture Notes in Computer Science*, pages 32–44, Turku Finland, July 2004. Springer-Verlag.  
<cs-www.cs.yale.edu/homes/aspnes/localization-abstract.html>

- [16] D. Avis and K. Fukuda. A pivoting algorithm for convex hulls and vertex enumeration of arrangements and polyhedra. *Discrete and Computational Geometry*, 8:295–313, 1992.  
<http://www.convexoptimization.com/TOOLS/AvisFukuda.pdf>
- [17] Christine Bachoc and Frank Vallentin. New upper bounds for kissing numbers from semidefinite programming. *Journal of the American Mathematical Society*, 21(3):909–924, July 2008.  
<http://arxiv.org/abs/math/0608426>
- [18] Mihály Bakonyi and Charles R. Johnson. The Euclidean distance matrix completion problem. *SIAM Journal on Matrix Analysis and Applications*, 16(2):646–654, April 1995.
- [19] Keith Ball. An elementary introduction to modern convex geometry. In Silvio Levy, editor, *Flavors of Geometry*, volume 31, chapter 1, pages 1–58. MSRI Publications, 1997.  
[www.msri.org/publications/books/Book31/files/ball.pdf](http://www.msri.org/publications/books/Book31/files/ball.pdf)
- [20] Richard G. Baraniuk. Compressive sensing [lecture notes]. *IEEE Signal Processing Magazine*, 24(4):118–121, July 2007.  
<http://www.convexoptimization.com/TOOLS/GeometryCardinality.pdf>
- [21] George Phillip Barker. Theory of cones. *Linear Algebra and its Applications*, 39:263–291, 1981.
- [22] George Phillip Barker and David Carlson. Cones of diagonally dominant matrices. *Pacific Journal of Mathematics*, 57(1):15–32, 1975.
- [23] George Phillip Barker and James Foran. Self-dual cones in Euclidean spaces. *Linear Algebra and its Applications*, 13:147–155, 1976.
- [24] Dror Baron, Michael B. Wakin, Marco F. Duarte, Shriram Sarvotham, and Richard G. Baraniuk. Distributed compressed sensing. Technical Report ECE-0612, Rice University, Electrical and Computer Engineering Department, December 2006.  
<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.85.4789>
- [25] Alexander I. Barvinok. Problems of distance geometry and convex properties of quadratic maps. *Discrete & Computational Geometry*, 13(2):189–202, 1995.
- [26] Alexander I. Barvinok. A remark on the rank of positive semidefinite matrices subject to affine constraints. *Discrete & Computational Geometry*, 25(1):23–31, 2001.  
<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.34.4627>
- [27] Alexander I. Barvinok. *A Course in Convexity*. American Mathematical Society, 2002.
- [28] Alexander I. Barvinok. Approximating orthogonal matrices by permutation matrices. *Pure and Applied Mathematics Quarterly*, 2:943–961, 2006.
- [29] Heinz H. Bauschke and Jonathan M. Borwein. On projection algorithms for solving convex feasibility problems. *SIAM Review*, 38(3):367–426, September 1996.
- [30] Steven R. Bell. *The Cauchy Transform, Potential Theory, and Conformal Mapping*. CRC Press, 1992.
- [31] Jean Bellissard and Bruno Iochum. Homogeneous and facially homogeneous self-dual cones. *Linear Algebra and its Applications*, 19:1–16, 1978.
- [32] Richard Bellman and Ky Fan. On systems of linear inequalities in Hermitian matrix variables. In Victor L. Klee, editor, *Convexity*, volume VII of *Proceedings of Symposia in Pure Mathematics*, pages 1–11. American Mathematical Society, 1963.
- [33] Adi Ben-Israel. Linear equations and inequalities on finite dimensional, real or complex, vector spaces: A unified theory. *Journal of Mathematical Analysis and Applications*, 27:367–389, 1969.
- [34] Adi Ben-Israel. Motzkin’s transposition theorem, and the related theorems of Farkas, Gordan and Stiemke. In Michiel Hazewinkel, editor, *Encyclopaedia of Mathematics*. Springer-Verlag, 2001.  
<http://www.convexoptimization.com/TOOLS/MOTZKIN.pdf>
- [35] Aharon Ben-Tal and Arkadi Nemirovski. *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*. SIAM, 2001.
- [36] Aharon Ben-Tal and Arkadi Nemirovski. Non-Euclidean restricted memory level method for large-scale convex optimization. *Mathematical Programming*, 102(3):407–456, January 2005.  
[http://www2.isye.gatech.edu/~nemirovs/Bundle-Mirror\\_rev.fin.pdf](http://www2.isye.gatech.edu/~nemirovs/Bundle-Mirror_rev.fin.pdf)
- [37] John J. Benedetto and Paulo J.S.G. Ferreira editors. *Modern Sampling Theory: Mathematics and Applications*. Birkhäuser, 2001.
- [38] Christian R. Berger, Javier Areta, Krishna Pattipati, and Peter Willett. Compressed sensing – A look beyond linear programming. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 3857–3860, 2008.

- [39] Radu Berinde, Anna C. Gilbert, Piotr Indyk, Howard J. Karloff, and Martin J. Strauss. Combining geometry and combinatorics: A unified approach to sparse signal recovery. In *46<sup>th</sup> Annual Allerton Conference on Communication, Control, and Computing*, pages 798–805. IEEE, September 2008.  
<http://arxiv.org/abs/0804.4666>
- [40] Abraham Berman. *Cones, Matrices, and Mathematical Programming*, volume 79 of *Lecture Notes in Economics and Mathematical Systems*. Springer-Verlag, 1973.
- [41] Abraham Berman and Naomi Shaked-Monderer. *Completely Positive Matrices*. World Scientific, 2003.
- [42] Dimitri P. Bertsekas. *Nonlinear Programming*. Athena Scientific, second edition, 1999.
- [43] Dimitri P. Bertsekas, Angelia Nedić, and Asuman E. Ozdaglar. *Convex Analysis and Optimization*. Athena Scientific, 2003.
- [44] Rajendra Bhatia. *Matrix Analysis*. Springer-Verlag, 1997.
- [45] Zhengbing Bian, Fabian Chudak, Robert Israel, Brad Lackey, William G. Macready, and Aidan Roy. Discrete optimization using quantum annealing on sparse Ising models. *frontiers in Physics*, 2:1–10, September 2014.  
<http://www.convexoptimization.com/TOOLS/frontiers.pdf>
- [46] Zhengbing Bian, Fabian Chudak, William G. Macready, and Geordie Rose. The Ising model: teaching an old problem new tricks. *D:Wave Systems*, August 2010.  
<http://www.dwavesys.com/sites/default/files/weightedmaxsat.v2.pdf>
- [47] Pratik Biswas, Tzu-Chen Liang, Kim-Chuan Toh, Yinyu Ye, and Ta-Chung Wang. Semidefinite programming approaches for sensor network localization with noisy distance measurements. *IEEE Transactions on Automation Science and Engineering*, 3(4):360–371, October 2006.
- [48] Pratik Biswas, Tzu-Chen Liang, Ta-Chung Wang, and Yinyu Ye. Semidefinite programming based algorithms for sensor network localization. *ACM Transactions on Sensor Networks*, 2(2):188–220, May 2006.  
[http://web.stanford.edu/~yyye/combined\\_rev3.pdf](http://web.stanford.edu/~yyye/combined_rev3.pdf)
- [49] Pratik Biswas, Kim-Chuan Toh, and Yinyu Ye. A distributed SDP approach for large-scale noisy anchor-free graph realization with applications to molecular conformation. *SIAM Journal on Scientific Computing*, 30(3):1251–1277, March 2008.  
<http://web.stanford.edu/~yyye/SISC-molecule-revised-2.pdf>
- [50] Pratik Biswas and Yinyu Ye. Semidefinite programming for ad hoc wireless sensor network localization. In *Proceedings of the Third International Symposium on Information Processing in Sensor Networks (IPSN)*, pages 46–54. IEEE, April 2004.  
<http://web.stanford.edu/~yyye/adhocn4.pdf>
- [51] Sudhendu Biswas. *Textbook of Matrix Algebra*. PHI, 3<sup>rd</sup> edition, 2012.  
<http://books.google.com/books?id=FxoA6Q2UJKwC>  
<http://www.convexoptimization.com/TOOLS/Vitulli.pdf>
- [52] Åke Björck and Tommy Elfving. Accelerated projection methods for computing pseudoinverse solutions of systems of linear equations. *BIT Numerical Mathematics*, 19(2):145–163, June 1979.  
<http://www.convexoptimization.com/TOOLS/BIT.pdf>
- [53] Leonard M. Blumenthal. A note on the four-point property. *Bulletin of the American Mathematical Society*, 39(6):423–426, 1933.  
<http://www.convexoptimization.com/TOOLS/euclid.pdf>
- [54] Leonard M. Blumenthal. *Theory and Applications of Distance Geometry*. Oxford University Press, 1953.
- [55] A. W. Bojanczyk and A. Lutoborski. The Procrustes problem for orthogonal Stiefel matrices. *SIAM Journal on Scientific Computing*, 21(4):1291–1304, December 1999.  
<http://citesearx.ist.psu.edu/viewdoc/summary?doi=10.1.1.15.9063>
- [56] Ingwer Borg and Patrick Groenen. *Modern Multidimensional Scaling*. Springer-Verlag, 1997.
- [57] Jonathan M. Borwein and Heinz Bauschke. Projection algorithms and monotone operators, 1998.  
<http://www.cecm.sfu.ca/~jborwein/projections4.pdf>
- [58] Jonathan M. Borwein and Adrian S. Lewis. *Convex Analysis and Nonlinear Optimization: Theory and Examples*. Springer-Verlag, 2000.  
<http://www.convexoptimization.com/TOOLS/Borwein.pdf>
- [59] Jonathan M. Borwein and Warren B. Moors. Stability of closedness of convex cones under linear mappings. *Journal of Convex Analysis*, 16(3):699–705, 2009.  
<http://www.convexoptimization.com/TOOLS/BorweinMoors.pdf>

- [60] Radu Ioan Boț, Ernö Robert Csetnek, and Gert Wanka. Regularity conditions via quasi-relative interior in convex programming. *SIAM Journal on Optimization*, 19(1):217–233, 2008.  
<http://www.convexoptimization.com/TOOLS/Wanka.pdf>
- [61] Richard Bouldin. The pseudo-inverse of a product. *SIAM Journal on Applied Mathematics*, 24(4):489–495, June 1973.  
<http://www.convexoptimization.com/TOOLS/Pseudoinverse.pdf>
- [62] Stephen Boyd and Jon Dattorro. Alternating projections. Stanford University, 2003.  
[http://web.stanford.edu/class/ee392o/alt\\_proj.pdf](http://web.stanford.edu/class/ee392o/alt_proj.pdf)
- [63] Stephen Boyd, Laurent El Ghaoui, Eric Feron, and Venkataramanan Balakrishnan. *Linear Matrix Inequalities in System and Control Theory*. SIAM, 1994.  
<http://web.stanford.edu/~boyd/lmibook>
- [64] Stephen Boyd, Seung-Jean Kim, Lieven Vandenberghe, and Arash Hassibi. A tutorial on geometric programming. *Optimization and Engineering*, 8(1):67–127, March 2007.  
[http://web.stanford.edu/~boyd/papers/gp\\_tutorial.html](http://web.stanford.edu/~boyd/papers/gp_tutorial.html)
- [65] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.  
<http://web.stanford.edu/~boyd/cvxbook>
- [66] James P. Boyle and Richard L. Dykstra. A method for finding projections onto the intersection of convex sets in Hilbert spaces. In R. Dykstra, T. Robertson, and F. T. Wright, editors, *Advances in Order Restricted Statistical Inference*, pages 28–47. Springer-Verlag, 1986.
- [67] Lev M. Brègman. The method of successive projection for finding a common point of convex sets. *Soviet Mathematics*, 162(3):487–490, 1965. AMS translation of Doklady Akademii Nauk SSSR, 6:688–692.
- [68] Lev M. Brègman, Yair Censor, Simeon Reich, and Yael Zepkowitz-Malachi. Finding the projection of a point onto the intersection of convex sets via projections onto halfspaces. *Journal of Approximation Theory*, 124(2):194–218, October 2003.  
[http://www.optimization-online.org/DB\\_FILE/2003/06/669.pdf](http://www.optimization-online.org/DB_FILE/2003/06/669.pdf)
- [69] Mike Brookes. Matrix reference manual: Matrix calculus, 2002.  
<http://www.ee.ic.ac.uk/hp/staff/dmb/matrix/intro.html>
- [70] J. P. Brooks, J. H. Dulá, and E. L. Boone. A pure  $L_1$ -norm principal component analysis. *Computational Statistics and Data Analysis*, 61:83–98, May 2013.  
<http://www.convexoptimization.com/TOOLS/Boone.pdf>
- [71] Richard A. Brualdi. *Combinatorial Matrix Classes*. Cambridge University Press, 2006.
- [72] Richard C. Cabot and Bruce E. Hofer *et alii*. *AA 501 Distortion Analyzer*. Tektronix, 1980.  
<http://www.convexoptimization.com/TOOLS/AA501.pdf>
- [73] Jian-Feng Cai, Emmanuel J. Candès, and Zuowei Shen. A singular value thresholding algorithm for matrix completion, October 2008.  
<http://arxiv.org/abs/0810.3286>
- [74] Alberto Cambini and Laura Martein. *Generalized Convexity and Optimization*. Springer-Verlag, 2009.
- [75] Emmanuel J. Candès. Compressive sampling. In *Proceedings of the International Congress of Mathematicians*, volume III, pages 1433–1452, Madrid, August 2006. European Mathematical Society.  
<http://www.acm.caltech.edu/~emmanuel/papers/CompressiveSampling.pdf>
- [76] Emmanuel J. Candès and Justin K. Romberg.  $\ell_1$ -MAGIC : Recovery of sparse signals via convex programming, October 2005.  
<https://statweb.stanford.edu/~candes/l1magic>
- [77] Emmanuel J. Candès, Justin K. Romberg, and Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, February 2006.  
<http://arxiv.org/abs/math.NA/0409186> (2004 DRAFT)  
<http://www.acm.caltech.edu/~emmanuel/papers/ExactRecovery.pdf>
- [78] Emmanuel J. Candès, Michael B. Wakin, and Stephen P. Boyd. Enhancing sparsity by reweighted  $\ell_1$  minimization. *Journal of Fourier Analysis and Applications*, 14(5-6):877–905, October 2008.  
<http://arxiv.org/abs/0711.1612v1>
- [79] Michael Carter, Holly Hui Jin, Michael Saunders, and Yinyu Ye. SPASELOC: An adaptive subproblem algorithm for scalable wireless sensor network localization. *SIAM Journal on Optimization*, 17(4):1102–1128, December 2006.  
<http://www.convexoptimization.com/TOOLS/JinSIAM.pdf>  
 Erratum: p.252 herein.

- [80] Lawrence Cayton and Sanjoy Dasgupta. Robust Euclidean embedding. In *Proceedings of the 23<sup>rd</sup> International Conference on Machine Learning* (ICML), Pittsburgh Pennsylvania USA, 2006.  
<http://cseweb.ucsd.edu/~lcayton/robEmb.pdf>
- [81] Yves Chabilliac and Jean-Pierre Crouzeix. Definiteness and semidefiniteness of quadratic forms revisited. *Linear Algebra and its Applications*, 63:283–292, 1984.
- [82] Manmohan K. Chandraker, Sameer Agarwal, Fredrik Kahl, David Nistér, and David J. Kriegman. Autocalibration via rank-constrained estimation of the absolute quadric. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, June 2007.  
<http://www2.maths.lth.se/vision/publdb/reports/pdf/chandraker-agarwal-etal-cvpr-07.pdf>
- [83] Rick Chartrand. Exact reconstruction of sparse signals via nonconvex minimization. *IEEE Signal Processing Letters*, 14(10):707–710, October 2007.  
<math.lanl.gov/Research/Publications/Docs/chartrand-2007-exact.pdf>
- [84] Chi-Tsong Chen. *Linear System Theory and Design*. Oxford University Press, 1999.
- [85] Scott ShaoBing Chen, David L. Donoho, and Michael A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1998.
- [86] Ward Cheney and Allen A. Goldstein. Proximity maps for convex sets. *Proceedings of the American Mathematical Society*, 10:448–450, 1959. Erratum: p.602 herein.
- [87] Mung Chiang. Geometric programming for communication systems. *Communications and Information Theory*, 2(1&2):1–154, July 2005.  
<https://www.princeton.edu/~chiangm/gp.pdf>
- [88] Stéphane Chrétien. An alternating l1 approach to the compressed sensing problem, September 2009.  
[http://arxiv.org/PS\\_cache/arxiv/pdf/0809/0809.0660v3.pdf](http://arxiv.org/PS_cache/arxiv/pdf/0809/0809.0660v3.pdf)
- [89] Steven Chu. Autobiography from *Les Prix Nobel*, 1997.  
[nobelprize.org/nobel\\_prizes/physics/lauriates/1997/chu-autobio.html](nobelprize.org/nobel_prizes/physics/lauriates/1997/chu-autobio.html)
- [90] Ruel V. Churchill and James Ward Brown. *Complex Variables and Applications*. McGraw-Hill, fifth edition, 1990.
- [91] Jon F. Claerbout and Francis Muir. Robust modeling of erratic data. *Geophysics*, 38:826–844, 1973.
- [92] J. H. Conway and N. J. A. Sloane. *Sphere Packings, Lattices and Groups*. Springer-Verlag, third edition, 1999.
- [93] John B. Conway. *A Course in Functional Analysis*. Springer-Verlag, second edition, 1990.
- [94] Jose A. Costa, Neal Patwari, and Alfred O. Hero III. Distributed weighted-multidimensional scaling for node localization in sensor networks. *ACM Transactions on Sensor Networks*, 2(1):39–64, February 2006.  
<http://www.convexoptimization.com/TOOLS/Costa.pdf>
- [95] Richard W. Cottle, Jong-Shi Pang, and Richard E. Stone. *The Linear Complementarity Problem*. Academic Press, 1992.
- [96] G. M. Crippen and T. F. Havel. *Distance Geometry and Molecular Conformation*. Wiley, 1988.
- [97] Frank Critchley. *Multidimensional scaling: a critical examination and some new proposals*. PhD thesis, University of Oxford, Nuffield College, 1980.
- [98] Frank Critchley. On certain linear mappings between inner-product and squared-distance matrices. *Linear Algebra and its Applications*, 105:91–107, 1988.  
<http://www.convexoptimization.com/TOOLS/Critchley.pdf>
- [99] Ronald E. Crochiere and Lawrence R. Rabiner. *Multirate Digital Signal Processing*. Prentice-Hall, 1983.
- [100] Lawrence B. Crowell. *Quantum Fluctuations of Spacetime*. World Scientific, 2005.
- [101] E. D. Dahl. Programming with D-Wave: Map coloring problem. *D:Wave Systems*, November 2013.  
<http://www.dwavesys.com/sites/default/files/Map%20Coloring%20WP2.pdf>
- [102] Joachim Dahl, Bernard H. Fleury, and Lieven Vandenberghe. Approximate maximum-likelihood estimation using semidefinite programming. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume VI, pages 721–724, April 2003.  
<http://www.convexoptimization.com/TOOLS/Dahl.pdf>
- [103] George B. Dantzig. *Linear Programming and Extensions*. Princeton University Press, 1963 (Rand Corporation).
- [104] Alexandre d'Aspremont, Laurent El Ghaoui, Michael I. Jordan, and Gert R. G. Lanckriet. A direct formulation for sparse PCA using semidefinite programming. In Lawrence K. Saul, Yair Weiss, and Léon Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 41–48. MIT Press, Cambridge Massachusetts USA, 2005.  
[http://books.nips.cc/papers/files/nips17/NIPS2004\\_0645.pdf](http://books.nips.cc/papers/files/nips17/NIPS2004_0645.pdf)

- [105] Jon Dattorro. Electronic apparatus for very high quality modification of the time duration and/or pitch of stereo music and speech. Lexicon Inc, 1985.  
<https://ccrma.stanford.edu/~dattorro/Lexipatent.htm>
- [106] Jon Dattorro. C program for *Constrained Least Squares Fit of a Filter Bank to an Arbitrary Magnitude Frequency Response*. Ensoniq Corp, 1991.  
<http://ccrma.stanford.edu/~dattorro/PhiLS.pdf>
- [107] Jon Dattorro. Constrained Least Squares Fit of a Filter Bank to an Arbitrary Magnitude Frequency Response. Ensoniq Corp, 1991.  
<http://ccrma.stanford.edu/~dattorro/Hearing.htm>
- [108] Jon Dattorro. Effect Design, Part 1: Reverberator and Other Filters. *Journal of the Audio Engineering Society*, 45(9):660–684, September 1997.  
<http://www.convexoptimization.com/TOOLS/EffectDesignPart1.pdf>
- [109] Jon Dattorro. Effect Design, Part 2: Delay-Line Modulation and Chorus. *Journal of the Audio Engineering Society*, 45(10):764–788, October 1997.  
<http://www.convexoptimization.com/TOOLS/EffectDesignPart2.pdf>
- [110] Jon Dattorro. Convex optimization of a first-order sigma delta modulator. Stanford University, 1999.  
<https://ccrma.stanford.edu/~dattorro/convex.pdf>
- [111] Jon Dattorro. Effect Design, Part 3: Oscillators: Sinusoidal and Pseudonoise. *Journal of the Audio Engineering Society*, 50(3):115–146, March 2002.  
<http://www.convexoptimization.com/TOOLS/EffectDesignPart3.pdf>
- [112] Jon Dattorro. Equality relating Euclidean distance cone to positive semidefinite cone. *Linear Algebra and its Applications*, 428(11+12):2597–2600, June 2008.  
<http://www.convexoptimization.com/TOOLS/DattorroLAA.pdf>
- [113] Joel Dawson, Stephen Boyd, Mar Hershenson, and Thomas Lee. Optimal allocation of local feedback in multistage amplifiers via geometric programming. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, 48(1):1–11, January 2001.  
[http://web.stanford.edu/~boyd/papers/fdbk\\_alloc.html](http://web.stanford.edu/~boyd/papers/fdbk_alloc.html)
- [114] Etienne de Klerk. *Aspects of Semidefinite Programming: Interior Point Algorithms and Selected Applications*. Kluwer Academic Publishers, 2002.
- [115] Jan de Leeuw. Fitting distances by least squares. UCLA Statistics Series Technical Report No. 130, Interdivisional Program in Statistics, UCLA, Los Angeles California USA, 1993.  
<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.50.7472>
- [116] Jan de Leeuw. Multidimensional scaling. In *International Encyclopedia of the Social & Behavioral Sciences*, pages 13512–13519. Elsevier, 2004.  
<http://www.convexoptimization.com/TOOLS/leeuw274.pdf>
- [117] Jan de Leeuw. Unidimensional scaling. In Brian S. Everitt and David C. Howell, editors, *Encyclopedia of Statistics in Behavioral Science*, volume 4, pages 2095–2097. Wiley, 2005.  
<http://www.convexoptimization.com/TOOLS/deLeeuw2005.pdf>
- [118] Jan de Leeuw and Willem Heiser. Theory of multidimensional scaling. In P. R. Krishnaiah and L. N. Kanal, editors, *Handbook of Statistics*, volume 2, chapter 13, pages 285–316. North-Holland Publishing, Amsterdam, 1982.
- [119] Erik D. Demaine and Martin L. Demaine. Jigsaw puzzles, edge matching, and polyomino packing: Connections and complexity. *Graphs and Combinatorics*, 23(supplement 1):195–208, Springer-Verlag, Tokyo, June 2007.
- [120] Erik D. Demaine, Francisco Gomez-Martin, Henk Meijer, David Rappaport, Perouz Taslakian, Godfried T. Toussaint, Terry Winograd, and David R. Wood. The distance geometry of music. *Computational Geometry: Theory and Applications*, 42(5):429–454, July 2009.  
<http://www.convexoptimization.com/TOOLS/dgm.pdf>
- [121] John R. D'Errico. DERIVEST, 2007.  
<http://www.convexoptimization.com/TOOLS/DERIVEST.pdf>
- [122] Frank Deutsch. *Best Approximation in Inner Product Spaces*. Springer-Verlag, 2001.
- [123] Frank Deutsch and Hein Hundal. The rate of convergence of Dykstra's cyclic projections algorithm: The polyhedral case. *Numerical Functional Analysis and Optimization*, 15:537–565, 1994.
- [124] Frank Deutsch and Peter H. Maserick. Applications of the Hahn-Banach theorem in approximation theory. *SIAM Review*, 9(3):516–530, July 1967.

- [125] Frank Deutsch, John H. McCabe, and George M. Phillips. Some algorithms for computing best approximations from convex cones. *SIAM Journal on Numerical Analysis*, 12(3):390–403, June 1975.
- [126] Michel Marie Deza and Monique Laurent. *Geometry of Cuts and Metrics*. Springer-Verlag, 1997. <http://www.convexoptimization.com/TOOLS/cutbook.pdf>
- [127] Sudhakar Waman Dharmadhikari and Kumar Joag-Dev. *Unimodality, Convexity, and Applications*. Academic Press, 1988.
- [128] Carolyn Pillers Dobler. A matrix approach to finding a set of generators and finding the polar (dual) of a class of polyhedral cones. *SIAM Journal on Matrix Analysis and Applications*, 15(3):796–803, July 1994. Erratum: p.164 herein.
- [129] Ivan Dokmanić. *Listening to Distances and Hearing Shapes: Inverse Problems in Room Acoustics and Beyond*. PhD thesis, École Polytechnique Fédérale De Lausanne, Laboratoire De Communications Audiovisuelles, 2015. <http://www.convexoptimization.com/TOOLS/Dokmanic.pdf>
- [130] Ivan Dokmanić, Reza Parhizkar, Juri Ranieri, and Martin Vetterli. Euclidean distance matrices: A short walk through theory, algorithms and applications, February 2015. <http://www.convexoptimization.com/TOOLS/EDMapplications.pdf>
- [131] Ivan Dokmanić, Reza Parhizkar, Juri Ranieri, and Martin Vetterli. Euclidean Distance Matrices: Essential theory, algorithms, and applications. *IEEE Signal Processing Magazine*, 32(6):12–30, November 2015. <http://www.convexoptimization.com/TOOLS/07298562.pdf>
- [132] Ivan Dokmanić, Reza Parhizkar, Andreas Walther, Yue M. Lu, and Martin Vetterli. Acoustic echoes reveal room shape. *Proceedings of the National Academy of Sciences*, 110(30):12186–12191, July 2013. <http://www.convexoptimization.com/TOOLS/pnas.pdf>
- [133] Elizabeth D. Dolan, Robert Fourer, Jorge J. Moré, and Todd S. Munson. Optimization on the NEOS server. *SIAM News*, 35(6):4,8,9, August 2002.
- [134] Bruce Randall Donald. 3-D structure in chemistry and molecular biology, 1998. <http://www.cs.duke.edu/brd/Teaching/Previous/Bio>
- [135] David L. Donoho. Neighborly polytopes and sparse solution of underdetermined linear equations. Technical Report 2005-04, Stanford University, Department of Statistics, January 2005. [www-stat.stanford.edu/~donoho/Reports/2005/NPaSSULE-01-28-05.pdf](http://www-stat.stanford.edu/~donoho/Reports/2005/NPaSSULE-01-28-05.pdf)
- [136] David L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, April 2006. [www-stat.stanford.edu/~donoho/Reports/2004/CompressedSensing091604.pdf](http://www-stat.stanford.edu/~donoho/Reports/2004/CompressedSensing091604.pdf)
- [137] David L. Donoho. High-dimensional centrally symmetric polytopes with neighborliness proportional to dimension. *Discrete & Computational Geometry*, 35(4):617–652, May 2006.
- [138] David L. Donoho, Michael Elad, and Vladimir Temlyakov. Stable recovery of sparse overcomplete representations in the presence of noise. *IEEE Transactions on Information Theory*, 52(1):6–18, January 2006. [www-stat.stanford.edu/~donoho/Reports/2004/StableSparse-Donoho-etal.pdf](http://www-stat.stanford.edu/~donoho/Reports/2004/StableSparse-Donoho-etal.pdf)
- [139] David L. Donoho and Philip B. Stark. Uncertainty principles and signal recovery. *SIAM Journal on Applied Mathematics*, 49(3):906–931, June 1989.
- [140] David L. Donoho and Jared Tanner. Neighborliness of randomly projected simplices in high dimensions. *Proceedings of the National Academy of Sciences*, 102(27):9452–9457, July 2005.
- [141] David L. Donoho and Jared Tanner. Sparse nonnegative solution of underdetermined linear equations by linear programming. *Proceedings of the National Academy of Sciences*, 102(27):9446–9451, July 2005. [www-stat.stanford.edu/~donoho/Reports/2004/SparseNonnegative-Donoho-Tanner.pdf](http://www-stat.stanford.edu/~donoho/Reports/2004/SparseNonnegative-Donoho-Tanner.pdf)
- [142] David L. Donoho and Jared Tanner. Counting faces of randomly projected polytopes when the projection radically lowers dimension. *Journal of the American Mathematical Society*, 22(1):1–53, January 2009.
- [143] Miguel Nuno Ferreira Fialho dos Anjos. *New Convex Relaxations for the Maximum Cut and VLSI Layout Problems*. PhD thesis, University of Waterloo, Ontario Canada, Department of Combinatorics and Optimization, 2001. <http://etd.uwaterloo.ca/etd/manjos2001.pdf>
- [144] John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra. Efficient projections onto the  $\ell_1$ -ball for learning in high dimensions. In *Proceedings of the 25<sup>th</sup> International Conference on Machine Learning* (ICML), pages 272–279, Helsinki Finland, July 2008. Association for Computing Machinery (ACM). <http://icml2008.cs.helsinki.fi/papers/361.pdf>

- [145] Richard L. Dykstra. An algorithm for restricted least squares regression. *Journal of the American Statistical Association*, 78(384):837–842, 1983.
- [146] Peter J. Eccles. *An Introduction to Mathematical Reasoning: numbers, sets and functions*. Cambridge University Press, 1997.
- [147] Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, September 1936.  
[www.convexoptimization.com/TOOLS/eckart&young.1936.pdf](http://www.convexoptimization.com/TOOLS/eckart&young.1936.pdf)
- [148] Alan Edelman, Tomás A. Arias, and Steven T. Smith. The geometry of algorithms with orthogonality constraints. *SIAM Journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.
- [149] John J. Edgell. Graphics calculator applications on 4-D constructs, 1996.  
<http://archives.math.utk.edu/ICTCM/EP-9/C47/pdf/paper.pdf>
- [150] Ivar Ekeland and Roger Témam. *Convex Analysis and Variational Problems*. SIAM, 1999.
- [151] Julius Farkas. Theorie der einfachen Ungleichungen. *Journal für die reine und angewandte Mathematik*, 124:1–27, 1902.  
[http://gdz.sub.uni-goettingen.de/dms/load/img/?PID=GDZPPN002165023&physid=PHYS\\_0006](http://gdz.sub.uni-goettingen.de/dms/load/img/?PID=GDZPPN002165023&physid=PHYS_0006)
- [152] Maryam Fazel. *Matrix Rank Minimization with Applications*. PhD thesis, Stanford University, Department of Electrical Engineering, March 2002.  
<http://faculty.washington.edu/mfazel/thesis-final.pdf>
- [153] Maryam Fazel, Haitham Hindi, and Stephen P. Boyd. A rank minimization heuristic with application to minimum order system approximation. In *Proceedings of the American Control Conference*, volume 6, pages 4734–4739. American Automatic Control Council (AACC), June 2001.  
<http://web.stanford.edu/~boyd/papers/nucnorm.html>
- [154] Maryam Fazel, Haitham Hindi, and Stephen P. Boyd. Log-det heuristic for matrix rank minimization with applications to Hankel and Euclidean distance matrices. In *Proceedings of the American Control Conference*, volume 3, pages 2156–2162. American Automatic Control Council (AACC), June 2003.  
[http://faculty.washington.edu/mfazel/acc03\\_final.pdf](http://faculty.washington.edu/mfazel/acc03_final.pdf)
- [155] Maryam Fazel, Haitham Hindi, and Stephen P. Boyd. Rank minimization and applications in system theory. In *Proceedings of the American Control Conference*, volume 4, pages 3273–3278. American Automatic Control Council (AACC), June 2004.  
<http://faculty.washington.edu/mfazel/acc04-tutorial.pdf>
- [156] J. T. Feddema, R. H. Byrne, J. J. Harrington, D. M. Kilman, C. L. Lewis, R. D. Robinett, B. P. Van Leeuwen, and J. G. Young. Advanced mobile networking, sensing, and controls. Sandia Report SAND2005-1661, Sandia National Laboratories, Albuquerque New Mexico USA, March 2005.  
[www.prod.sandia.gov/cgi-bin/techlib/access-control.pl/2005/051661.pdf](http://prod.sandia.gov/cgi-bin/techlib/access-control.pl/2005/051661.pdf)
- [157] César Fernández, Thomas Szyperski, Thierry Bruyère, Paul Ramage, Egon Mössinger, and Kurt Wüthrich. NMR solution structure of the pathogenesis-related protein P14a. *Journal of Molecular Biology*, 266:576–593, 1997.
- [158] O. P. Ferreira and S. Z. Németh. Generalized projections onto convex sets. *Journal of Global Optimization*, 52:831–842, 2012.  
<http://www.convexoptimization.com/TOOLS/Ferreira.pdf>
- [159] Richard Phillips Feynman, Robert B. Leighton, and Matthew L. Sands. *The Feynman Lectures on Physics: Commemorative Issue*, volume I. Addison-Wesley, 1989.
- [160] Anthony V. Fiacco and Garth P. McCormick. *Nonlinear Programming: Sequential Unconstrained Minimization Techniques*. SIAM, 1990.
- [161] Önder Filiz and Aylin Yener. Rank constrained temporal-spatial filters for CDMA systems with base station antenna arrays. In *Proceedings of the Johns Hopkins University Conference on Information Sciences and Systems*, March 2003.  
[http://wcan.ee.psu.edu/papers/filiz-yener\\_ciiss03.pdf](http://wcan.ee.psu.edu/papers/filiz-yener_ciiss03.pdf)
- [162] P. A. Fillmore and J. P. Williams. Some convexity theorems for matrices. *Glasgow Mathematical Journal*, 12:110–117, 1971.
- [163] Paul Finsler. Über das Vorkommen definiter und semidefiniter Formen in Scharen quadratischer Formen. *Commentarii Mathematici Helvetici*, 9:188–192, 1937.
- [164] P. E. Fleischer and J. Tow. Design formulas for biquad active filters using three operational amplifiers. *Proceedings of the IEEE*, 61(5):662–663, May 1973.  
<http://www.convexoptimization.com/TOOLS/FleischerTow.pdf>
- [165] Anders Forsgren, Philip E. Gill, and Margaret H. Wright. Interior methods for nonlinear optimization. *SIAM Review*, 44(4):525–597, 2002.

- [166] Shmuel Friedland and Anatoli Torokhti. Generalized rank-constrained matrix approximations. *SIAM Journal on Matrix Analysis and Applications*, 29(2):656–659, March 2007.  
<http://arxiv.org/pdf/math.0C/0603674.pdf>
- [167] Ragnar Frisch. The multiplex method for linear programming. *SANKHYĀ: The Indian Journal of Statistics*, 18(3/4):329–362, September 1957.  
<http://www.convexoptimization.com/TOOLS/Frisch57.pdf>
- [168] Norbert Gaffke and Rudolf Mathar. A cyclic projection algorithm via duality. *Metrika*, 36:29–54, 1989.
- [169] Jérôme Galtier. Semi-definite programming as a simple extension to linear programming: convex optimization with queueing, equity and other telecom functionals. In *3ème Rencontres Francophones sur les Aspects Algorithmiques des Télécommunications* (AlgoTel), pages 21–28. INRIA, Saint Jean de Luz FRANCE, May 2001.  
<http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.8.7378>
- [170] Bernd Gärtner and Jiří Matoušek. *Approximation Algorithms and Semidefinite Programming*. Springer-Verlag Berlin Heidelberg, 2012.
- [171] Laurent El Ghaoui. EE 227A: Convex Optimization and Applications, Lecture 11 – October 3. University of California, Berkeley, Fall 2006. Scribe: Nikhil Shetty.  
<http://www.convexoptimization.com/TOOLS/Ghaoui.pdf>
- [172] Laurent El Ghaoui and Silviu-Iulian Niculescu, editors. *Advances in Linear Matrix Inequality Methods in Control*. SIAM, 2000.
- [173] Philip E. Gill, Walter Murray, and Margaret H. Wright. *Practical Optimization*. Academic Press, 1981.
- [174] Philip E. Gill, Walter Murray, and Margaret H. Wright. *Numerical Linear Algebra and Optimization*, volume 1. Addison-Wesley, 1991.
- [175] James Gleik. *Isaac Newton*. Pantheon Books, 2003.
- [176] W. Glunt, Tom L. Hayden, S. Hong, and J. Wells. An alternating projection algorithm for computing the nearest Euclidean distance matrix. *SIAM Journal on Matrix Analysis and Applications*, 11(4):589–600, 1990.
- [177] K. Goebel and W. A. Kirk. *Topics in Metric Fixed Point Theory*. Cambridge University Press, 1990.
- [178] Michel X. Goemans and David P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the Association for Computing Machinery*, 42(6):1115–1145, November 1995.  
<http://www.convexoptimization.com/TOOLS/maxcut-jacm.pdf>
- [179] D. Goldfarb and K. Scheinberg. Interior point trajectories in semidefinite programming. *SIAM Journal on Optimization*, 8(4):871–886, 1998.
- [180] Gene Golub and William Kahan. Calculating the singular values and pseudo-inverse of a matrix. *SIAM Journal on Numerical Analysis*, Series B, 2(2):205–224, 1965.  
<http://www.convexoptimization.com/TOOLS/GolubKahan.pdf>
- [181] Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. Johns Hopkins, third edition, 1996.
- [182] Gene H. Golub and Urs von Matt. Tikhonov regularization for large scale problems, 1997.  
<http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.51.409>
- [183] P. Gordan. Ueber die auflösung linearer gleichungen mit reellen coefficienten. *Mathematische Annalen*, 6:23–28, 1873.
- [184] Irina F. Gorodnitsky and Bhaskar D. Rao. Sparse signal reconstruction from limited data using FOCUSS: A re-weighted minimum norm algorithm. *IEEE Transactions on Signal Processing*, 45(3):600–616, March 1997.  
<http://www.convexoptimization.com/TOOLS/focuss.pdf>
- [185] John Clifford Gower. Adding a point to vector diagrams in multivariate analysis. *Biometrika*, 55(3):582–585, 1968.  
<http://www.convexoptimization.com/TOOLS/GOWER1968.pdf>
- [186] John Clifford Gower. Euclidean distance geometry. *The Mathematical Scientist*, 7:1–14, 1982.  
<http://www.convexoptimization.com/TOOLS/Gower2.pdf>
- [187] John Clifford Gower. Properties of Euclidean and non-Euclidean distance matrices. *Linear Algebra and its Applications*, 67:81–97, 1985.  
<http://www.convexoptimization.com/TOOLS/Gower1.pdf>

- [188] John Clifford Gower and Garnt B. Dijksterhuis. *Procrustes Problems*. Oxford University Press, 2004.
- [189] John Clifford Gower and David J. Hand. *Biplots*. Chapman & Hall, 1996.
- [190] Alexander Graham. *Kronecker Products and Matrix Calculus with Applications*. Ellis Horwood Limited, 1981.
- [191] Michael Grant and Stephen Boyd.  
`cvx`: MATLAB software for disciplined convex programming, 2015.  
<http://cvxr.com>
- [192] Robert M. Gray. Toeplitz and circulant matrices: A review. *Foundations and Trends in Communications and Information Theory*, 2(3):155–239, 2006.  
<http://www-ee.stanford.edu/~gray/toeplitz.pdf>
- [193] T. N. E. Greville. Note on the generalized inverse of a matrix product. *SIAM Review*, 8:518–521, 1966.
- [194] Rémi Gribonval and Morten Nielsen. Sparse representations in unions of bases. *IEEE Transactions on Information Theory*, 49(12):3320–3325, December 2003.  
[http://people.math.aau.dk/~mnielsen/reprints/sparse\\_unions.pdf](http://people.math.aau.dk/~mnielsen/reprints/sparse_unions.pdf)
- [195] Rémi Gribonval and Morten Nielsen. Highly sparse representations from dictionaries are unique and independent of the sparseness measure. *Applied and Computational Harmonic Analysis*, 22(3):335–355, May 2007.  
<http://www.convexoptimization.com/TOOLS/R-2003-16.pdf>
- [196] Karolos M. Grigoriadis and Eric B. Beran. Alternating projection algorithms for linear matrix inequalities problems with rank constraints. In Laurent El Ghaoui and Silviu-Iulian Niculescu, editors, *Advances in Linear Matrix Inequality Methods in Control*, chapter 13, pages 251–267. SIAM, 2000.
- [197] Peter Gritzmann and Victor Klee. On the complexity of some basic problems in computational convexity: II. Volume and mixed volumes. Technical Report TR:94-31, DIMACS, Rutgers University, 1994.  
<http://dimacs.rutgers.edu/TechnicalReports/TechReports/1994/94-31.ps>
- [198] Peter Gritzmann and Victor Klee. On the complexity of some basic problems in computational convexity: II. Volume and mixed volumes. In T. Bisztriczky, P. McMullen, R. Schneider, and A. Ivić Weiss, editors, *Polytopes: Abstract, Convex and Computational*, pages 373–466. Kluwer Academic Publishers, 1994.
- [199] L. G. Gubin, B. T. Polyak, and E. V. Raik. The method of projections for finding the common point of convex sets. *U.S.S.R. Computational Mathematics and Mathematical Physics*, 7(6):1–24, 1967.
- [200] Osman Güler and Yinyu Ye. Convergence behavior of interior-point algorithms. *Mathematical Programming*, 60(2):215–228, 1993.
- [201] P. R. Halmos. Positive approximants of operators. *Indiana University Mathematics Journal*, 21:951–960, 1972.
- [202] Shih-Ping Han. A successive projection method. *Mathematical Programming*, 40:1–14, 1988.
- [203] Godfrey H. Hardy, John E. Littlewood, and George Pólya. *Inequalities*. Cambridge University Press, second edition, 1952.
- [204] R. Harris and M. W. Johnson *et alii*. Experimental investigation of an eight-qubit unit cell in a superconducting optimization processor. *Physical Review B - Condensed Matter and Materials Physics*, 82(2-024511), July 2010.  
<http://www.convexoptimization.com/TOOLS/Harris.pdf>
- [205] Mahir Hassan and Amir Khajepour. Layout and force optimisation in cable-driven parallel manipulators. In Lihui Wang and Jeff Xi, editors, *Smart Devices and Machines for Advanced Manufacturing*, pages 110–135. Springer-Verlag, 2010.  
[https://books.google.com/books?id=GRxmIwN9u\\_UC&printsec=frontcover](https://books.google.com/books?id=GRxmIwN9u_UC&printsec=frontcover)
- [206] Arash Hassibi and Mar Hershenson. Automated optimal design of switched-capacitor filters. In *Proceedings of the Conference on Design, Automation, and Test in Europe*, page 1111, March 2002.
- [207] Johan Hästads. Some optimal inapproximability results. In *Proceedings of the Twenty-Ninth Annual ACM Symposium on Theory of Computing* (STOC), pages 1–10, El Paso Texas USA, 1997. Association for Computing Machinery (ACM).  
<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.16.183>, 2002.
- [208] Johan Hästads. Some optimal inapproximability results. *Journal of the Association for Computing Machinery*, 48(4):798–859, July 2001.

- [209] Timothy F. Havel and Kurt Wüthrich. An evaluation of the combined use of nuclear magnetic resonance and distance geometry for the determination of protein conformations in solution. *Journal of Molecular Biology*, 182:281–294, 1985.
- [210] Tom L. Hayden and Jim Wells. Approximation by matrices positive semidefinite on a subspace. *Linear Algebra and its Applications*, 109:115–130, 1988.
- [211] Tom L. Hayden, Jim Wells, Wei-Min Liu, and Pablo Tarazaga. The cone of distance matrices. *Linear Algebra and its Applications*, 144:153–169, 1991.
- [212] Uwe Helmke and John B. Moore. *Optimization and Dynamical Systems*. Springer-Verlag, 1994.
- [213] Bruce Hendrickson. Conditions for unique graph realizations. *SIAM Journal on Computing*, 21(1):65–84, February 1992.
- [214] T. Herrmann, Peter Güntert, and Kurt Wüthrich. Protein NMR structure determination with automated NOE assignment using the new software CANDID and the torsion angle dynamics algorithm DYANA. *Journal of Molecular Biology*, 319(1):209–227, May 2002.
- [215] Mar Hershenson. Design of pipeline analog-to-digital converters via geometric programming. In *Proceedings of the IEEE/ACM International Conference on Computer Aided Design (ICCAD)*, pages 317–324, November 2002.
- [216] Mar Hershenson. Efficient description of the design space of analog circuits. In *Proceedings of the 40<sup>th</sup> ACM/IEEE Design Automation Conference*, pages 970–973, June 2003.
- [217] Mar Hershenson, Stephen Boyd, and Thomas Lee. Optimal design of a CMOS OpAmp via geometric programming. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 20(1):1–21, January 2001.  
<http://web.stanford.edu/~boyd/papers/opamp.html>
- [218] Mar Hershenson, Dave Colleran, Arash Hassibi, and Navraj Nandra. Synthesizable full custom mixed-signal IP. *Electronics Design Automation Consortium (EDA)*, 2002.  
<http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.520.7705>
- [219] Mar Hershenson, Sunderarajan S. Mohan, Stephen Boyd, and Thomas Lee. Optimization of inductor circuits via geometric programming. In *Proceedings of the 36<sup>th</sup> ACM/IEEE Design Automation Conference*, pages 994–998, June 1999.  
[http://web.stanford.edu/~boyd/papers/inductor\\_opt.html](http://web.stanford.edu/~boyd/papers/inductor_opt.html)
- [220] Nick Higham. Matrix Procrustes problems, 1995.  
<http://www.convexoptimization.com/TOOLS/procrust94.ps>  
 Lecture notes.
- [221] Richard D. Hill and Steven R. Waters. On the cone of positive semidefinite matrices. *Linear Algebra and its Applications*, 90:81–88, 1987.
- [222] Jean-Baptiste Hiriart-Urruty. Ensembles de Tchebychev vs. ensembles convexes: l'état de la situation vu via l'analyse convexe non lisse. *Annales des Sciences Mathématiques du Québec*, 22(1):47–62, 1998.
- [223] Jean-Baptiste Hiriart-Urruty. Global optimality conditions in maximizing a convex quadratic function under convex quadratic constraints. *Journal of Global Optimization*, 21(4):445–455, December 2001.  
[http://www.convexoptimization.com/TOOLS/Jean\\_Baptiste.pdf](http://www.convexoptimization.com/TOOLS/Jean_Baptiste.pdf)
- [224] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Convex Analysis and Minimization Algorithms II: Advanced Theory and Bundle Methods*. Springer-Verlag, second edition, 1996.
- [225] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of Convex Analysis*. Springer-Verlag, 2001.
- [226] Alan J. Hoffman and Helmut W. Wielandt. The variation of the spectrum of a normal matrix. *Duke Mathematical Journal*, 20:37–40, 1953.
- [227] Alfred Horn. Doubly stochastic matrices and the diagonal of a rotation matrix. *American Journal of Mathematics*, 76(3):620–630, July 1954.  
<http://www.convexoptimization.com/TOOLS/AHorn.pdf>
- [228] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 1987.
- [229] Roger A. Horn and Charles R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, 1994.
- [230] Alston S. Householder. *The Theory of Matrices in Numerical Analysis*. Dover, 1975.

- [231] Yifan Hu. Adjacency matrix graph corresponding to smallest face of Eternity II problem, 2011.  
<http://www.convexoptimization.com/TOOLS/EternityIImovie.eps>  
<http://www.convexoptimization.com/TOOLS/EternityIImovie.pdf>
- [232] Hong-Xuan Huang, Zhi-An Liang, and Panos M. Pardalos. Some properties for the Euclidean distance matrix and positive semidefinite matrix completion problems. *Journal of Global Optimization*, 25(1):3–21, January 2003.  
<http://www.convexoptimization.com/TOOLS/pardalos.pdf>
- [233] Lawrence Hubert, Jacqueline Meulman, and Willem Heiser. Two purposes for matrix factorization: A historical appraisal. *SIAM Review*, 42(1):68–82, 2000.  
<http://www.convexoptimization.com/TOOLS/hubert.pdf>
- [234] Xiaoming Huo. *Sparse Image Representation via Combined Transforms*. PhD thesis, Stanford University, Department of Statistics, August 1999.  
<http://www.convexoptimization.com/TOOLS/ReweightingFrom1999XiaomingHuo.pdf>
- [235] 5W Infographic. Wireless 911. *Technology Review*, 107(5):78–79, June 2004.  
<http://www.technologyreview.com/communications/13643>
- [236] George Isac. *Complementarity Problems*. Springer-Verlag, 1992.
- [237] Matt Jacobson. Analyze n-dimensional polyhedra in terms of vertices or (in)equalities, 2015.  
<https://www.mathworks.com/matlabcentral/fileexchange/30892>
- [238] Nathan Jacobson. *Lectures in Abstract Algebra, vol.II - Linear Algebra*. Van Nostrand, 1953.
- [239] Viren Jain and Lawrence K. Saul. Exploratory analysis and visualization of speech and music by locally linear embedding. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 3, pages 984–987, May 2004.  
[http://www.cs.ucsd.edu/~saul/papers/lle\\_icassp04.pdf](http://www.cs.ucsd.edu/~saul/papers/lle_icassp04.pdf)
- [240] Joakim Jaldén. Bi-criterion  $\ell_1/\ell_2$ -norm optimization. Master's thesis, Royal Institute of Technology (KTH), Department of Signals Sensors and Systems, Stockholm Sweden, September 2002.  
<http://www.convexoptimization.com/TOOLS/JaldenMSThesis.pdf>
- [241] Joakim Jaldén, Cristoff Martin, and Björn Ottersten. Semidefinite programming for detection in linear systems - Optimality conditions and space-time decoding. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume IV, pages 9–12, April 2003.  
<http://www.convexoptimization.com/TOOLS/Jalden.pdf>
- [242] Florian Jarre. Convex analysis on symmetric matrices. In Henry Wolkowicz, Romesh Saigal, and Lieven Vandenberghe, editors, *Handbook of Semidefinite Programming: Theory, Algorithms, and Applications*, chapter 2. Kluwer, 2000.  
<http://www.convexoptimization.com/TOOLS/Handbook.pdf>
- [243] Holly Hui Jin. *Scalable Sensor Localization Algorithms for Wireless Sensor Networks*. PhD thesis, University of Toronto, Graduate Department of Mechanical and Industrial Engineering, 2005.  
[web.stanford.edu/group/SOL/dissertations/holly-thesis.pdf](http://web.stanford.edu/group/SOL/dissertations/holly-thesis.pdf)
- [244] Charles R. Johnson and Pablo Tarazaga. Connections between the real positive semidefinite and distance matrix completion problems. *Linear Algebra and its Applications*, 223/224:375–391, 1995.
- [245] Charles R. Johnson and Pablo Tarazaga. Binary representation of normalized symmetric and correlation matrices. *Linear and Multilinear Algebra*, 52(5):359–366, 2004.
- [246] M. W. Johnson and M. H. S. Amin *et alii*. Quantum annealing with manufactured spins. *Nature*, 473:194–198, May 2011.  
<http://www.convexoptimization.com/TOOLS/manufacturedspins.pdf>
- [247] Mark Kahrs and Karlheinz Brandenburg, editors. *Applications of Digital Signal Processing to Audio and Acoustics*. Kluwer Academic Publishers, 1998.
- [248] Thomas Kailath. *Linear Systems*. Prentice-Hall, 1980.
- [249] Tosio Kato. *Perturbation Theory for Linear Operators*. Springer-Verlag, 1966.
- [250] Paul J. Kelly and Norman E. Ladd. *Geometry*. Scott, Foresman and Company, 1965.
- [251] Sunyoung Kim, Masakazu Kojima, Hayato Waki, and Makoto Yamashita. SFSDP: a Sparse version of Full SemiDefinite Programming relaxation for sensor network localization problems. Research Report B-457, Department of Mathematical and Computing Sciences, Tokyo Institute of Technology, Japan, July 2009.  
<http://math.ewha.ac.kr/~skim/Research/B-457.pdf>

- [252] Yoonsoo Kim and Mehran Mesbahi. On the rank minimization problem. In *Proceedings of the American Control Conference*, volume 3, pages 2015–2020. American Automatic Control Council (AACC), June 2004.
- [253] Ron Kimmel. *Numerical Geometry of Images: Theory, Algorithms, and Applications*. Springer-Verlag, 2003.
- [254] Erwin Kreyszig. *Introductory Functional Analysis with Applications*. Wiley, 1989.
- [255] M. J. Kronenburg. The binomial coefficient for negative arguments, March 2015.  
<https://arxiv.org/pdf/1105.3689.pdf>
- [256] Anthony Kuh, Chaopin Zhu, and Danilo Mandic. Sensor network localization using least squares kernel regression. In 10<sup>th</sup> International Conference of Knowledge-Based Intelligent Information and Engineering Systems (KES), volume 4253(III) of *Lecture Notes in Computer Science*, pages 1280–1287, Bournemouth UK, October 2006. Springer-Verlag.  
<http://www.convexoptimization.com/TOOLS/Kuh.pdf>
- [257] Harold W. Kuhn. Nonlinear programming: a historical view. In Richard W. Cottle and Carlton E. Lemke, editors, *Nonlinear Programming*, pages 1–26. American Mathematical Society, 1976.
- [258] Takahito Kuno, Yasutoshi Yajima, and Hiroshi Konno. An outer approximation method for minimizing the product of several convex functions on a convex set. *Journal of Global Optimization*, 3(3):325–335, September 1993.
- [259] Amy N. Langville and Carl D. Meyer. *Google’s PageRank and Beyond: The Science of Search Engine Rankings*. Princeton University Press, 2006.
- [260] Jean B. Lasserre. A new Farkas lemma for positive semidefinite matrices. *IEEE Transactions on Automatic Control*, 40(6):1131–1133, June 1995.
- [261] Jean B. Lasserre and Eduardo S. Zeron. A Laplace transform algorithm for the volume of a convex polytope. *Journal of the Association for Computing Machinery*, 48(6):1126–1140, November 2001.  
<http://arxiv.org/abs/math/0106168> (title revision).
- [262] Alan J. Laub. *Matrix Analysis for Scientists and Engineers*. SIAM, 2005.
- [263] Monique Laurent. A connection between positive semidefinite and Euclidean distance matrix completion problems. *Linear Algebra and its Applications*, 273:9–22, 1998.
- [264] Monique Laurent. A tour d’horizon on positive semidefinite and Euclidean distance matrix completion problems. In Panos M. Pardalos and Henry Wolkowicz, editors, *Topics in Semidefinite and Interior-Point Methods*, pages 51–76. American Mathematical Society, 1998.
- [265] Monique Laurent. Matrix completion problems. In Christodoulos A. Floudas and Panos M. Pardalos, editors, *Encyclopedia of Optimization*, volume III(Interior-M), pages 221–229. Kluwer, 2001.  
<http://homepages.cwi.nl/~monique/files/opt.ps>
- [266] Monique Laurent and Svatopluk Poljak. On a positive semidefinite relaxation of the cut polytope. *Linear Algebra and its Applications*, 223/224:439–461, 1995.
- [267] Monique Laurent and Svatopluk Poljak. On the facial structure of the set of correlation matrices. *SIAM Journal on Matrix Analysis and Applications*, 17(3):530–547, July 1996.
- [268] Monique Laurent and Franz Rendl. Semidefinite programming and integer programming. *Optimization Online*, 2002.  
[http://www.optimization-online.org/DB\\_HTML/2002/12/585.html](http://www.optimization-online.org/DB_HTML/2002/12/585.html)
- [269] Monique Laurent and Franz Rendl. Semidefinite programming and integer programming. In K. Aardal, George L. Nemhauser, and R. Weismantel, editors, *Discrete Optimization*, volume 12 of *Handbooks in Operations Research and Management Science*, chapter 8, pages 393–514. Elsevier, 2005.
- [270] Charles L. Lawson and Richard J. Hanson. *Solving Least Squares Problems*. SIAM, 1995.
- [271] Jung Rye Lee. The law of cosines in a tetrahedron. *Journal of the Korea Society of Mathematical Education, Series B (The Pure and Applied Mathematics)*, 4(1):1–6, 1997.
- [272] Claude Lemaréchal. Note on an extension of “Davidon” methods to nondifferentiable functions. *Mathematical Programming*, 7(1):384–387, December 1974.  
<http://www.convexoptimization.com/TOOLS/Lemarechal.pdf>
- [273] Vladimir L. Levin. Quasi-convex functions and quasi-monotone operators. *Journal of Convex Analysis*, 2(1/2):167–172, 1995.
- [274] Scott Nathan Levine. *Audio Representations for Data Compression and Compressed Domain Processing*. PhD thesis, Stanford University, Department of Electrical Engineering, 1999.  
[http://www.convexoptimization.com/TOOLS/Levine\\_Thesis.pdf](http://www.convexoptimization.com/TOOLS/Levine_Thesis.pdf)

- [275] Doron Levy. Introduction to numerical analysis, September 2010.  
<http://www.math.umd.edu/~dlevy/books/na.pdf>
- [276] Adrian S. Lewis. Eigenvalue-constrained faces. *Linear Algebra and its Applications*, 269:159–181, 1998.
- [277] Anhua Lin. *Projection algorithms in nonlinear programming*. PhD thesis, Johns Hopkins University, 2003.
- [278] Miguel Sousa Lobo, Lieven Vandenberghe, Stephen Boyd, and Hervé Lebret. Applications of second-order cone programming. *Linear Algebra and its Applications*, 284:193–228, November 1998. Special Issue on Linear Algebra in Control, Signals and Image Processing.  
<http://web.stanford.edu/~boyd/papers/socp.html>
- [279] Lee Lorch and Donald J. Newman. On the composition of completely monotonic functions and completely monotonic sequences and related questions. *Journal of the London Mathematical Society*, second series, 28:31–45, 1983.  
<http://www.convexoptimization.com/TOOLS/Lorch.pdf>
- [280] David G. Luenberger. *Optimization by Vector Space Methods*. Wiley, 1969.
- [281] David G. Luenberger. *Introduction to Dynamic Systems: Theory, Models, & Applications*. Wiley, 1979.
- [282] David G. Luenberger. *Linear and Nonlinear Programming*. Addison-Wesley, second edition, 1989.
- [283] Zhi-Quan Luo, Jos F. Sturm, and Shuzhong Zhang. Superlinear convergence of a symmetric primal-dual path following algorithm for semidefinite programming. *SIAM Journal on Optimization*, 8(1):59–81, 1998.
- [284] Zhi-Quan Luo and Wei Yu. An introduction to convex optimization for communications and signal processing. *IEEE Journal On Selected Areas In Communications*, 24(8):1426–1438, August 2006.
- [285] Morris Marden. *Geometry of Polynomials*. American Mathematical Society, second edition, 1985.
- [286] K. V. Mardia. Some properties of classical multi-dimensional scaling. *Communications in Statistics: Theory and Methods*, A7(13):1233–1241, 1978.  
<http://www.convexoptimization.com/TOOLS/CommunStatTheo1978.pdf>
- [287] K. V. Mardia, J. T. Kent, and J. M. Bibby. *Multivariate Analysis*. Academic Press, 1979.
- [288] Robert J. Marks II, editor. *Advanced Topics in Shannon Sampling and Interpolation Theory*. Springer-Verlag, 1993.
- [289] Jerrold E. Marsden and Michael J. Hoffman. *Elementary Classical Analysis*. Freeman, second edition, 1995.
- [290] Rudolf Mathar. The best Euclidean fit to a given distance matrix in prescribed dimensions. *Linear Algebra and its Applications*, 67:1–6, 1985.
- [291] Rudolf Mathar. *Multidimensionale Skalierung*. B. G. Teubner Stuttgart, 1997.
- [292] Nathan S. Mendelsohn and A. Lloyd Dulmage. The convex hull of sub-permutation matrices. *Proceedings of the American Mathematical Society*, 9(2):253–254, April 1958.  
<http://www.convexoptimization.com/TOOLS/permu.pdf>
- [293] Mehran Mesbahi and G. P. Papavassiliopoulos. On the rank minimization problem over a positive semi-definite linear matrix inequality. *IEEE Transactions on Automatic Control*, 42(2):239–243, February 1997.
- [294] Mehran Mesbahi and G. P. Papavassiliopoulos. Solving a class of rank minimization problems via semi-definite programs, with applications to the fixed order output feedback synthesis. In *Proceedings of the American Control Conference*, volume 1, pages 77–80. American Automatic Control Council (AACC), June 1997.  
<http://www.convexoptimization.com/TOOLS/Mesbahi.pdf>
- [295] Sunderarajan S. Mohan, Mar Hershenson, Stephen Boyd, and Thomas Lee. Simple accurate expressions for planar spiral inductances. *IEEE Journal of Solid-State Circuits*, 34(10):1419–1424, October 1999.  
[http://web.stanford.edu/~boyd/papers/inductance\\_expressions.html](http://web.stanford.edu/~boyd/papers/inductance_expressions.html)
- [296] Sunderarajan S. Mohan, Mar Hershenson, Stephen Boyd, and Thomas Lee. Bandwidth extension in CMOS with optimized on-chip inductors. *IEEE Journal of Solid-State Circuits*, 35(3):346–355, March 2000.  
[http://web.stanford.edu/~boyd/papers/bandwidth\\_ext.html](http://web.stanford.edu/~boyd/papers/bandwidth_ext.html)

- [297] David Moore, John Leonard, Daniela Rus, and Seth Teller. Robust distributed network localization with noisy range measurements. In *Proceedings of the Second International Conference on Embedded Networked Sensor Systems* (SenSys'04), pages 50–61, Baltimore Maryland USA, November 2004. Association for Computing Machinery (ACM, Winner of the Best Paper Award). <http://rvsn.csail.mit.edu/netloc/sensys04.pdf>
- [298] E. H. Moore. On the reciprocal of the general algebraic matrix. *Bulletin of the American Mathematical Society*, 26:394–395, 1920. Abstract.
- [299] B. S. Mordukhovich. Maximum principle in the problem of time optimal response with nonsmooth constraints. *Journal of Applied Mathematics and Mechanics*, 40:960–969, 1976.
- [300] Jean-Jacques Moreau. Décomposition orthogonale d'un espace Hilbertien selon deux cônes mutuellement polaires. *Comptes Rendus de l'Académie des Sciences, Paris*, 255:238–240, 1962.
- [301] T. S. Motzkin and I. J. Schoenberg. The relaxation method for linear inequalities. *Canadian Journal of Mathematics*, 6:393–404, 1954.
- [302] Neil Muller, Lourenço Magaia, and B. M. Herbst. Singular value decomposition, eigenfaces, and 3D reconstructions. *SIAM Review*, 46(3):518–545, September 2004.
- [303] Katta G. Murty and Feng-Tien Yu. *Linear Complementarity, Linear and Nonlinear Programming*. Heldermann Verlag, Internet edition, 1988. [http://www-personal.umich.edu/~murty/books/linear\\_complementarity\\_webbook](http://www-personal.umich.edu/~murty/books/linear_complementarity_webbook)
- [304] Oleg R. Musin. An extension of Delsarte's method. The kissing problem in three and four dimensions. In *Proceedings of the 2004 COE Workshop on Sphere Packings*, Kyushu University Japan, 2005. <http://arxiv.org/abs/math.MG/0512649>
- [305] Stephen G. Nash and Ariela Sofer. *Linear and Nonlinear Programming*. McGraw-Hill, 1996.
- [306] John Lawrence Nazareth. *Differentiable Optimization and Equation Solving: A Treatise on Algorithmic Science and the Karmarkar Revolution*. Springer-Verlag, 2003.
- [307] A. B. Németh and S. Z. Németh. How to project onto the monotone nonnegative cone using pool adjacent violators type algorithms, January 2012. <http://arxiv.org/abs/1201.2343>
- [308] Arkadi Nemirovski. *Lectures on Modern Convex Optimization*. <http://www.convexoptimization.com/TOOLS/LectModConvOpt.pdf>, 2005.
- [309] Yurii Nesterov and Arkadii Nemirovskii. *Interior-Point Polynomial Algorithms in Convex Programming*. SIAM, 1994.
- [310] Ewa Niewiadomska-Szynkiewicz and Michał Marks. Optimization schemes for wireless sensor network localization. *International Journal of Applied Mathematics and Computer Science*, 19(2):291–302, 2009. <http://www.convexoptimization.com/TOOLS/Marks.pdf>
- [311] L. Nirenberg. *Functional Analysis*. New York University, New York, 1961. Lectures given in 1960–1961, notes by Lesley Sibner.
- [312] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer-Verlag, 1999.
- [313] Royal Swedish Academy of Sciences. Nobel prize in chemistry, 2002. [www.nobelprize.org/nobel\\_prizes/chemistry/laureates/2002/popular.html](http://www.nobelprize.org/nobel_prizes/chemistry/laureates/2002/popular.html)
- [314] C. S. Ogilvy. *Excursions in Geometry*. Dover, 1990. Citation: *Proceedings of the CUPM Geometry Conference*, Mathematical Association of America, No.16 (1967), p.21.
- [315] Ricardo Oliveira, João Costeira, and João Xavier. Optimal point correspondence through the use of rank constraints. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, pages 1016–1021, June 2005.
- [316] Alan V. Oppenheim and Ronald W. Schafer. *Discrete-Time Signal Processing*. Prentice-Hall, 1989.
- [317] Robert Orsi, Uwe Helmke, and John B. Moore. A Newton-like method for solving rank constrained linear matrix inequalities. *Automatica*, 42(11):1875–1882, 2006. <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.384.5385>
- [318] Brad Osgood. Notes on the Ahlfors mapping of a multiply connected domain, 2000. [www-ee.stanford.edu/~osgood/Ahlfors-Bergman-Szeg.pdf](http://www-ee.stanford.edu/~osgood/Ahlfors-Bergman-Szeg.pdf)
- [319] M. L. Overton and R. S. Womersley. On the sum of the largest eigenvalues of a symmetric matrix. *SIAM Journal on Matrix Analysis and Applications*, 13:41–45, 1992.
- [320] Pythagoras Papadimitriou. *Parallel Solution of SVD-Related Problems, With Applications*. PhD thesis, Department of Mathematics, University of Manchester, October 1993.

- [321] Panos M. Pardalos and Henry Wolkowicz, editors. *Topics in Semidefinite and Interior-Point Methods*. American Mathematical Society, 1998.
- [322] Reza Parhizkar. *Euclidean Distance Matrices: Properties, Algorithms and Applications*. PhD thesis, École Polytechnique Fédérale De Lausanne, Laboratoire De Communications Audiovisuelles, 2013. <http://www.convexoptimization.com/TOOLS/Parhizkar.pdf>
- [323] Bradford W. Parkinson, James J. Spilker, Penina Axelrad, and Per Enge, editors. *Global Positioning System: Theory and Applications, Volume I*. American Institute of Aeronautics and Astronautics, 1996.
- [324] Beresford Parlett and Gilbert Strang. Matrices with prescribed Ritz values. *Linear Algebra and its Applications*, 428(7):1725–1739, April 2008. <http://www.convexoptimization.com/TOOLS/ArrowheadMatrixStrang.pdf>
- [325] Gábor Pataki. Cone-LP's and semidefinite programs: Geometry and a simplex-type method. In William H. Cunningham, S. Thomas McCormick, and Maurice Queyranne, editors, *Proceedings of the 5<sup>th</sup> International Conference on Integer Programming and Combinatorial Optimization* (IPCO), volume 1084 of *Lecture Notes in Computer Science*, pages 162–174, Vancouver, British Columbia, Canada, June 1996. Springer-Verlag.
- [326] Gábor Pataki. On the rank of extreme matrices in semidefinite programs and the multiplicity of optimal eigenvalues. *Mathematics of Operations Research*, 23(2):339–358, 1998. <http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.49.3831>  
Erratum: p.534 herein.
- [327] Gábor Pataki. The geometry of semidefinite programming. In Henry Wolkowicz, Romesh Saigal, and Lieven Vandenberghe, editors, *Handbook of Semidefinite Programming: Theory, Algorithms, and Applications*, chapter 3. Kluwer, 2000. <http://www.convexoptimization.com/TOOLS/Handbook.pdf>
- [328] Cengiz Pehlevan, Tao Hu, and Dmitri B. Chklovskii. A Hebbian/anti-Hebbian neural network for linear subspace learning: A derivation from multidimensional scaling of streaming data. *Neural Computation*, 27(7):1461–1495, July 2015. <https://www.researchgate.net/publication/273003026>
- [329] Teemu Pennanen and Jonathan Eckstein. Generalized Jacobians of vector-valued convex functions. Technical Report RRR 6-97, RUTCOR, Rutgers University, May 1997. <http://rutcor.rutgers.edu/pub/rrr/reports97/06.ps>
- [330] Roger Penrose. A generalized inverse for matrices. In *Proceedings of the Cambridge Philosophical Society*, volume 51, pages 406–413, 1955.
- [331] Chris Perkins. A convergence analysis of Dykstra's algorithm for polyhedral sets. *SIAM Journal on Numerical Analysis*, 40(2):792–804, 2002.
- [332] Sam Perlis. *Theory of Matrices*. Addison-Wesley, 1958.
- [333] Kaare Brandt Petersen and Michael Syskind Pedersen. *The Matrix Cookbook*, November 2012. <http://www.convexoptimization.com/TOOLS/matrixcookbook.pdf>
- [334] Florian Pfender and Günter M. Ziegler. Kissing numbers, sphere packings, and some unexpected proofs. *Notices of the American Mathematical Society*, 51(8):873–883, September 2004. <http://www.convexoptimization.com/TOOLS/Pfender.pdf>
- [335] Benjamin Recht. *Convex Modeling with Priors*. PhD thesis, Massachusetts Institute of Technology, Media Arts and Sciences Department, 2006. <http://www.convexoptimization.com/TOOLS/06.06.Recht.pdf>
- [336] Benjamin Recht, Maryam Fazel, and Pablo A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *Optimization Online*, June 2007. <http://faculty.washington.edu/mfazel/low-rank-v2.pdf>
- [337] Alex Reznik. Problem 3.41, from Sergio Verdú. *Multiuser Detection*. Cambridge University Press, 1998. [www.ee.princeton.edu/~verdu/mud/solutions/3/3.41.areznik.pdf](http://www.ee.princeton.edu/~verdu/mud/solutions/3/3.41.areznik.pdf), 2001.
- [338] Arthur Wayne Roberts and Dale E. Varberg. *Convex Functions*. Academic Press, 1973.
- [339] Sara Robinson. Hooked on meshing, researcher creates award-winning triangulation program. *SIAM News*, 36(9), November 2003. <http://www.siam.org/news/news.php?id=370>
- [340] R. Tyrrell Rockafellar. *Conjugate Duality and Optimization*. SIAM, 1974.
- [341] R. Tyrrell Rockafellar. Lagrange multipliers in optimization. *SIAM-American Mathematical Society Proceedings*, 9:145–168, 1976.

- [342] R. Tyrrell Rockafellar. Lagrange multipliers and optimality. *SIAM Review*, 35(2):183–238, June 1993.
- [343] R. Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, 1997 (first published in 1970). <http://www.convexoptimization.com/TOOLS/ConvexAnalysisRockafellar.pdf>
- [344] C. K. Rushforth. Signal restoration, functional analysis, and Fredholm integral equations of the first kind. In Henry Stark, editor, *Image Recovery: Theory and Application*, chapter 1, pages 1–27. Academic Press, 1987.
- [345] Peter Santos. Application platforms and synthesizable analog IP. In *Proceedings of the Embedded Systems Conference*, San Francisco California USA, 2003. <http://www.convexoptimization.com/TOOLS/Santos.pdf>
- [346] Shankar Sastry. *Nonlinear Systems: Analysis, Stability, and Control*. Springer-Verlag, 1999.
- [347] Uwe Schäfer. A linear complementarity problem with a P-matrix. *SIAM Review*, 46(2):189–201, June 2004.
- [348] Adriaan M. J. Schakel. Quantum phase transitions in 2d quantum liquids. In Diana V. Shopova and Dimo I. Uzunov, editors, *Correlations, Coherence, and Order*, pages 295–356. Kluwer Academic / Plenum Publishers, 1999.
- [349] Isaac J. Schoenberg. Remarks to Maurice Fréchet’s article “Sur la définition axiomatique d’une classe d’espaces distanciés vectoriellement applicable sur l’espace de Hilbert”. *Annals of Mathematics*, 36(3):724–732, July 1935. <http://www.convexoptimization.com/TOOLS/Schoenberg2.pdf>
- [350] Isaac J. Schoenberg. Metric spaces and positive definite functions. *Transactions of the American Mathematical Society*, 44:522–536, 1938. <http://www.convexoptimization.com/TOOLS/Schoenberg3.pdf>
- [351] Peter H. Schönemann. A generalized solution of the orthogonal Procrustes problem. *Psychometrika*, 31(1):1–10, March 1966. <http://www.convexoptimization.com/TOOLS/Schonemann.pdf>
- [352] Peter H. Schönemann, Tim Dorcey, and K. Kienapple. Subadditive concatenation in dissimilarity judgements. *Perception and Psychophysics*, 38:1–17, 1985.
- [353] Alexander Schrijver. On the history of combinatorial optimization (till 1960). In K. Aardal, G. L. Nemhauser, and R. Weismantel, editors, *Handbook of Discrete Optimization*, pages 1–68. Elsevier, 2005. <http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.22.3362>
- [354] Seymour Sherman. Doubly stochastic matrices and complex vector spaces. *American Journal of Mathematics*, 77(2):245–246, April 1955. <http://www.convexoptimization.com/TOOLS/Sherman.pdf>
- [355] Joshua A. Singer. *Log-Penalized Linear Regression*. PhD thesis, Stanford University, Department of Electrical Engineering, June 2004. <http://www.convexoptimization.com/TOOLS/Josh.pdf>
- [356] Anthony Man-Cho So and Yinyu Ye. A semidefinite programming approach to tensegrity theory and realizability of graphs. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithms* (SODA), pages 766–775, Miami Florida USA, January 2006. Association for Computing Machinery (ACM). <http://www.convexoptimization.com/TOOLS/YeSo.pdf>
- [357] Anthony Man-Cho So and Yinyu Ye. Theory of semidefinite programming for sensor network localization. *Mathematical Programming*, Series **A** and **B**, 109(2):367–384, January 2007. <http://web.stanford.edu/~yyye/local-theory.pdf>
- [358] D. C. Sorensen. Newton’s method with a model trust region modification. *SIAM Journal on Numerical Analysis*, 19(2):409–426, April 1982. <http://www.convexoptimization.com/TOOLS/soren.pdf>
- [359] Nathan Srebro, Jason D. M. Rennie, and Tommi S. Jaakkola. Maximum-margin matrix factorization. In Lawrence K. Saul, Yair Weiss, and Léon Bottou, editors, *Advances in Neural Information Processing Systems 17* (NIPS 2004), pages 1329–1336. MIT Press, 2005.
- [360] Wolfram Stadler, editor. *Multicriteria Optimization in Engineering and in the Sciences*. Springer-Verlag, 1988.
- [361] Henry Stark, editor. *Image Recovery: Theory and Application*. Academic Press, 1987.

- [362] Henry Stark. Polar, spiral, and generalized sampling and interpolation. In Robert J. Marks II, editor, *Advanced Topics in Shannon Sampling and Interpolation Theory*, chapter 6, pages 185–218. Springer-Verlag, 1993.
- [363] Willi-Hans Steeb. *Matrix Calculus and Kronecker Product with Applications and C++ Programs*. World Scientific, 1997.
- [364] Gilbert W. Stewart and Ji-guang Sun. *Matrix Perturbation Theory*. Academic Press, 1990.
- [365] Josef Stoer. On the characterization of least upper bound norms in matrix space. *Numerische Mathematik*, 6(1):302–314, December 1964.  
<http://www.convexoptimization.com/TOOLS/Stoer.pdf>
- [366] Josef Stoer and Christoph Witzgall. *Convexity and Optimization in Finite Dimensions I*. Springer-Verlag, 1970.
- [367] Nevena Jakovčević Stor, Ivan Slapničar, and Jesse L. Barlow. Accurate eigenvalue decomposition of real symmetric arrowhead matrices and applications. *Linear Algebra and its Applications*, 464:62–89, January 2015.  
<http://arxiv.org/abs/1302.7203>
- [368] Gilbert Strang. *Linear Algebra and its Applications*. Harcourt Brace, third edition, 1988.
- [369] Gilbert Strang. *Calculus*. Wellesley-Cambridge Press, 1992.
- [370] Gilbert Strang. *Introduction to Linear Algebra*. Wellesley-Cambridge Press, second edition, 1998.
- [371] Gilbert Strang. Course 18.06: Linear algebra, 2004.  
<http://web.mit.edu/18.06/www/Course-Info/Tcodes.html>
- [372] Stefan Straszewicz. Über exponierte Punkte abgeschlossener Punktmengen. *Fundamenta Mathematicae*, 24:139–143, 1935.  
<http://www.convexoptimization.com/TOOLS/Straszewicz.pdf>
- [373] Jos F. Sturm. SeDuMi (self-dual minimization). Software for optimization over symmetric cones, 2003.  
<http://sedumi.ie.lehigh.edu>
- [374] Jos F. Sturm and Shuzhong Zhang. On cones of nonnegative quadratic functions. *Mathematics of Operations Research*, 28(2):246–267, May 2003.  
[www.optimization-online.org/DB\\_HTML/2001/05/324.html](http://www.optimization-online.org/DB_HTML/2001/05/324.html)
- [375] George P. H. Styan. A review and some extensions of Takemura's generalizations of Cochran's theorem. Technical Report 56, Stanford University, Department of Statistics, September 1982.
- [376] George P. H. Styan and Akimichi Takemura. Rank additivity and matrix polynomials. Technical Report 57, Stanford University, Department of Statistics, September 1982.
- [377] Jun Sun, Stephen Boyd, Lin Xiao, and Persi Diaconis. The fastest mixing Markov process on a graph and a connection to a maximum variance unfolding problem. *SIAM Review*, 48(4):681–699, December 2006.  
<http://stanford.edu/~boyd/papers/fmmp.html>
- [378] Rangarajan K. Sundaram. *A First Course in Optimization Theory*. Cambridge University Press, 1996.
- [379] Chen Han Sung and Bit-Shun Tam. A study of projectionally exposed cones. *Linear Algebra and its Applications*, 139:225–252, 1990.
- [380] Yoshio Takane. On the relations among four methods of multidimensional scaling. *Behaviormetrika*, 4:29–43, 1977.  
<http://takane.brinkster.net/Yoshio/p008.pdf>
- [381] Akimichi Takemura. On generalizations of Cochran's theorem and projection matrices. Technical Report 44, Stanford University, Department of Statistics, August 1980.
- [382] Yasuhiko Takenaga and Toby Walsh. Tetravex is NP-complete. *Information Processing Letters*, 99(5):171–174, September 2006.  
<http://arxiv.org/pdf/0903.1147v1.pdf>
- [383] Dharmpal Takhar, Jason N. Laska, Michael B. Wakin, Marco F. Duarte, Dror Baron, Shriram Sarvotham, Kevin F. Kelly, and Richard G. Baraniuk. A new compressive imaging camera architecture using optical-domain compression. In *Proceedings of the SPIE Conference on Computational Imaging IV*, volume 6065, pages 43–52, February 2006.  
<http://www.convexoptimization.com/TOOLS/csCamera-SPIE-dec05.pdf>
- [384] Peng Hui Tan and Lars K. Rasmussen. The application of semidefinite programming for detection in CDMA. *IEEE Journal on Selected Areas in Communications*, 19(8), August 2001.

- [385] Pham Dinh Tao and Le Thi Hoai An. A D.C. optimization algorithm for solving the trust-region subproblem. *SIAM Journal on Optimization*, 8(2):476–505, 1998.  
<http://www.convexoptimization.com/TOOLS/pham.pdf>
- [386] Pablo Tarazaga. Faces of the cone of Euclidean distance matrices: Characterizations, structure and induced geometry. *Linear Algebra and its Applications*, 408:1–13, 2005.
- [387] George B. Thomas. *Calculus and Analytic Geometry*. Addison-Wesley, fourth edition, 1972.
- [388] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B (Methodological)*, 58(1):267–288, 1996.
- [389] Kim-Chuan Toh, Michael J. Todd, and Reha H. Tütüncü. SDPT3 – a MATLAB software for semidefinite-quadratic-linear programming, February 2009.  
<http://www.math.nus.edu.sg/~mattohkc/sdpt3.html>
- [390] Warren S. Torgerson. *Theory and Methods of Scaling*. Wiley, 1958.
- [391] Lloyd N. Trefethen and David Bau III. *Numerical Linear Algebra*. SIAM, 1997.
- [392] Michael W. Trosset. Applications of multidimensional scaling to molecular conformation. *Computing Science and Statistics*, 29:148–152, 1998.
- [393] Michael W. Trosset. Distance matrix completion by numerical optimization. *Computational Optimization and Applications*, 17(1):11–22, October 2000.
- [394] Michael W. Trosset. Extensions of classical multidimensional scaling: Computational theory. [convexoptimization.com/TOOLS/TrossetPITA.pdf](http://www.convexoptimization.com/TOOLS/TrossetPITA.pdf), 2001. Revision of technical report entitled “Computing distances between convex sets and subsets of the positive semidefinite matrices” first published in 1997.
- [395] Michael W. Trosset and Rudolf Mathar. On the existence of nonglobal minimizers of the STRESS criterion for metric multidimensional scaling. In *Proceedings of the Statistical Computing Section*, pages 158–162. American Statistical Association, 1997.
- [396] Joshua Trzasko and Armando Manduca. Highly undersampled magnetic resonance image reconstruction via homotopic  $\ell_0$ -minimization. *IEEE Transactions on Medical Imaging*, 28(1):106–121, January 2009.  
<http://www.convexoptimization.com/TOOLS/SparseRecon.pdf>
- [397] Joshua Trzasko, Armando Manduca, and Eric Borisch. Sparse MRI reconstruction via multiscale  $L_0$ -continuation. In *Proceedings of the IEEE 14<sup>th</sup> Workshop on Statistical Signal Processing*, pages 176–180, August 2007.  
<http://www.convexoptimization.com/TOOLS/L0MRI.pdf>
- [398] P. P. Vaidyanathan. *Multirate Systems and Filter Banks*. Prentice-Hall, 1993.
- [399] Jan van Tiel. *Convex Analysis, an Introductory Text*. Wiley, 1984.
- [400] Lieven Vandenberghe and Stephen Boyd. Semidefinite programming. *SIAM Review*, 38(1):49–95, March 1996.
- [401] Lieven Vandenberghe and Stephen Boyd. Connections between semi-infinite and semidefinite programming. In R. Reemtsen and J.-J. Rückmann, editors, *Semi-Infinite Programming*, chapter 8, pages 277–294. Kluwer Academic Publishers, 1998.
- [402] Lieven Vandenberghe and Stephen Boyd. Applications of semidefinite programming. *Applied Numerical Mathematics*, 29(3):283–299, March 1999.
- [403] Lieven Vandenberghe, Stephen Boyd, and Shao-Po Wu. Determinant maximization with linear matrix inequality constraints. *SIAM Journal on Matrix Analysis and Applications*, 19(2):499–533, April 1998.
- [404] Robert J. Vanderbei. Convex optimization: Interior-point methods and applications. Lecture notes, 1999.  
<http://orfe.princeton.edu/~rvdb/pdf/talks/pumath/talk.pdf>
- [405] Richard S. Varga. *Geršgorin and His Circles*. Springer-Verlag, 2004.
- [406] Martin Vetterli and Jelena Kovačević. *Wavelets and Subband Coding*. Prentice-Hall, 1995.
- [407] Martin Vetterli, Pina Marziliano, and Thierry Blu. Sampling signals with finite rate of innovation. *IEEE Transactions on Signal Processing*, 50(6):1417–1428, June 2002.  
<http://bigwww.epfl.ch/publications/vetterli0201.pdf>  
<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.19.5630>
- [408] È. B. Vinberg. The theory of convex homogeneous cones. *Transactions of the Moscow Mathematical Society*, 12:340–403, 1963. American Mathematical Society and London Mathematical Society joint translation, 1965.

- [409] John von Neumann. *Functional Operators, Volume II: The Geometry of Orthogonal Spaces*. Princeton University Press, 1950. Reprinted from mimeographed lecture notes first distributed in 1933.
- [410] *Wikimization*. An alternative proof without Moreau's theorem.  
[convexoptimization.com/wikimization/index.php/Complementarity\\_problem](http://convexoptimization.com/wikimization/index.php/Complementarity_problem)
- [411] *Wikimization*. Boolean feasibility.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [412] *Wikimization*. Compressive sampling of images by convex iteration.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [413] *Wikimization*. Conic independence.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [414] *Wikimization*. Convex iteration rank-1, 2013.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [415] *Wikimization*. Convex optimization of Eternity II.  
[convexoptimization.com/wikimization/index.php/Dattorro\\_Convex\\_Optimization\\_of\\_Eternity\\_II](http://convexoptimization.com/wikimization/index.php/Dattorro_Convex_Optimization_of_Eternity_II)
- [416] *Wikimization*. EDM using ordinal data.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [417] *Wikimization*. Fast max cut.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [418] *Wikimization*. Fast projection on monotone nonnegative cone.  
[convexoptimization.com/wikimization/index.php/Projection\\_on\\_Polyhedral\\_Cone](http://convexoptimization.com/wikimization/index.php/Projection_on_Polyhedral_Cone)
- [419] *Wikimization*. Fermat point.  
[http://www.convexoptimization.com/wikimization/index.php/Fermat\\_point](http://www.convexoptimization.com/wikimization/index.php/Fermat_point)
- [420] *Wikimization*. Filter design by convex iteration.  
[convexoptimization.com/wikimization/index.php/Filter\\_design\\_by\\_convex\\_iteration](http://convexoptimization.com/wikimization/index.php/Filter_design_by_convex_iteration)
- [421] *Wikimization*. Fixed point problems.  
[convexoptimization.com/wikimization/index.php/Complementarity\\_problem](http://convexoptimization.com/wikimization/index.php/Complementarity_problem)
- [422] *Wikimization*. High-order polynomials.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [423] *Wikimization*. Map of the USA.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [424] *Wikimization*. MATLAB for convex optimization & Euclidean distance geometry.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [425] *Wikimization*. Projection on simplicial cones.  
[convexoptimization.com/wikimization/index.php/Projection\\_on\\_Polyhedral\\_Cone](http://convexoptimization.com/wikimization/index.php/Projection_on_Polyhedral_Cone)
- [426] *Wikimization*. Rank reduction subroutine.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [427] *Wikimization*. Sampling sparsity.  
[convexoptimization.com/wikimization/index.php/Sampling\\_Sparsity](http://convexoptimization.com/wikimization/index.php/Sampling_Sparsity)
- [428] *Wikimization*. Signal dropout problem.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [429] *Wikimization*. Singular value decomposition (SVD) by rank-1 transformation.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [430] *Wikimization*. Sturm & Zhang's procedure for constructing dyad-decomposition.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [431] *Wikimization*. `isedm()`.  
[convexoptimization.com/wikimization/index.php/Matlab\\_for\\_Convex\\_Optimization](http://convexoptimization.com/wikimization/index.php/Matlab_for_Convex_Optimization)
- [432] Michael B. Wakin, Jason N. Laska, Marco F. Duarte, Dror Baron, Shriram Sarvotham, Dharmpal Takhar, Kevin F. Kelly, and Richard G. Baraniuk. An architecture for compressive imaging. In *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, pages 1273–1276, October 2006.  
<http://www.convexoptimization.com/TOOLS/CSCam-ICIP06.pdf>
- [433] Dominic Walliman. *How The Quantum Annealing Process Works*. D:Wave Systems, 2015.  
[https://www.youtube.com/watch?v=UV\\_RlCAC5Zs](https://www.youtube.com/watch?v=UV_RlCAC5Zs)
- [434] Roger Webster. *Convexity*. Oxford University Press, 1994.

- [435] Kilian Q. Weinberger and Lawrence K. Saul. Unsupervised learning of image manifolds by semidefinite programming. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, pages 988–995, 2004.  
[http://www.cs.ucsd.edu/~saul/papers/sde\\_cvpr04.pdf](http://www.cs.ucsd.edu/~saul/papers/sde_cvpr04.pdf)
- [436] Eric W. Weisstein. Mathworld – A Wolfram Web Resource, 2007.  
<http://mathworld.wolfram.com>
- [437] D. J. White. An analogue derivation of the dual of the general Fermat problem. *Management Science*, 23(1):92–94, September 1976.  
<http://www.convexoptimization.com/TOOLS/White.pdf>
- [438] Norbert Wiener. On factorization of matrices. *Commentarii Mathematici Helvetici*, 29:97–111, 1955.
- [439] Ami Wiesel, Yonina C. Eldar, and Shlomo Shamai (Shitz). Semidefinite relaxation for detection of 16-QAM signaling in MIMO channels. *IEEE Signal Processing Letters*, 12(9):653–656, September 2005.
- [440] Michael P. Williamson, Timothy F. Havel, and Kurt Wüthrich. Solution conformation of proteinase inhibitor IIA from bull seminal plasma by  $^1\text{H}$  nuclear magnetic resonance and distance geometry. *Journal of Molecular Biology*, 182:295–315, 1985.
- [441] Willie W. Wong. Cayley-Menger determinant and generalized N-dimensional Pythagorean theorem, November 2003. Application of Linear Algebra: Notes on Talk given to Princeton University Math Club.  
<http://www.convexoptimization.com/TOOLS/gp-r.pdf>
- [442] William Wooton, Edwin F. Beckenbach, and Frank J. Fleming. *Modern Analytic Geometry*. Houghton Mifflin, 1975.
- [443] Margaret H. Wright. The interior-point revolution in optimization: History, recent developments, and lasting consequences. *Bulletin of the American Mathematical Society*, 42(1):39–56, January 2005.
- [444] Stephen J. Wright. *Primal-Dual Interior-Point Methods*. SIAM, 1997.
- [445] Shao-Po Wu. *max-det Programming with Applications in Magnitude Filter Design*. A dissertation submitted to the department of Electrical Engineering, Stanford University, December 1997.
- [446] Shao-Po Wu and Stephen Boyd. `sdpSol`: A parser/solver for semidefinite programming and determinant maximization problems with matrix structure, May 1996.  
[http://web.stanford.edu/~boyd/old\\_software/SDPSOL.html](http://web.stanford.edu/~boyd/old_software/SDPSOL.html)
- [447] Shao-Po Wu and Stephen Boyd. `sdpSol`: A parser/solver for semidefinite programs with matrix structure. In Laurent El Ghaoui and Silviu-Iulian Niculescu, editors, *Advances in Linear Matrix Inequality Methods in Control*, chapter 4, pages 79–91. SIAM, 2000.  
<http://web.stanford.edu/~boyd/papers/sdpSol.html>
- [448] Shao-Po Wu, Stephen Boyd, and Lieven Vandenberghe. FIR filter design via spectral factorization and convex optimization. In Biswa Nath Datta, editor, *Applied and Computational Control, Signals, and Circuits*, volume 1, chapter 5, pages 215–246. Birkhäuser, 1998.  
<http://web.stanford.edu/~boyd/papers/magdes.html>
- [449] Naoki Yamamoto and Maryam Fazel. Computational approach to quantum encoder design for purity optimization. *Physical Review A*, 76(1), July 2007.  
<http://arxiv.org/abs/quant-ph/0606106>
- [450] David D. Yao, Shuzhong Zhang, and Xun Yu Zhou. Stochastic linear-quadratic control via primal-dual semidefinite programming. *SIAM Review*, 46(1):87–111, March 2004. Erratum: p.183 herein.
- [451] Yinyu Ye. Semidefinite programming for Euclidean distance geometric optimization. Lecture notes, 2003.  
<http://web.stanford.edu/class/ee392o/EE392o-yinyu-ye.pdf>
- [452] Yinyu Ye. Convergence behavior of central paths for convex homogeneous self-dual cones, 1996.  
<http://web.stanford.edu/~yyye/yyye/ye.ps>
- [453] Yinyu Ye. *Interior Point Algorithms: Theory and Analysis*. Wiley, 1997.
- [454] D. C. Youla. Mathematical theory of image restoration by the method of convex projection. In Henry Stark, editor, *Image Recovery: Theory and Application*, chapter 2, pages 29–77. Academic Press, 1987.
- [455] Fuzhen Zhang. *Matrix Theory: Basic Results and Techniques*. Springer-Verlag, 1999.

- [456] Kenneth M. Zick, Omar Shehab, and Matthew French. Experimental quantum annealing: case study involving the graph isomorphism problem. *Scientific Reports*, 5(11168), June 2015.  
<http://www.nature.com/articles/srep11168>
- [457] Günter M. Ziegler. Kissing numbers: Surprises in dimension four. *Emissary*, pages 4–5, Spring 2004.  
<http://www.msri.org/communications/emissary>

# Index

- $\emptyset$ , *see* set, empty
- $0$ , *see* zero, *see* origin
- $0$ -norm, 176, 179, 234, 271, 275, 280–283, 304, 635  
    Boolean, 234
- $1$ , 632
- $1$ -norm, 173–177, 179, 203, 249, 281, 472, 635  
    ball, *see* ball  
    Boolean, 234  
    distance matrix, 202, 203, 394  
    gradient, 179, 192, 559  
    heuristic, 272, 273, 472  
    nonnegative, 176, 177, 275, 280–284  
    projection, 604  
    signed, 285, 286
- $2$ -norm, 42, 170, 171, 173, 174, 286, 559, 635  
    ball, *see* ball  
    Boolean, 234  
    matrix, 58, 191, 249, 458, 468, 509, 516, 521, 530, 604, 635  
    constraint, 58  
    dual, 530, 635  
    inequality, 635  
    inverse, 530  
    Schur-form, 191, 530
- $\infty$ -norm, 173, 174, 179, 249, 635  
    ball, *see* ball  
    dual, 635
- $k$ -largest norm, 178, 179, 274, 321, 635  
    gradient, 180
- $\ell$ -norm, 635
- $\succ$ , 51, 81, 84, 91, 634
- $\succeq$ , 51, 84, 91, 634
- $\geq$ , 51, 635
- $\in$ , 623
- $\ni$ , 623
- $\perp$ , 42, 138, 628
- $\langle \rangle$ , 42, 43, 625
- $Ax = b$ , 41, 70, 174, 175, 179, 237, 238, 242, 246, 271, 274, 275, 332, 568–570  
    Boolean, 234  
    nonnegative, 71, 72, 176, 222, 275, 280–284  
    signed, 285, 286
- $\partial$ , 35, 39, 477, 543, 629
- $\delta$ , 485–487, 493, 495, 503, 569, 626
- $e$ , 48, 632
- $\varepsilon$ , 623
- $\exists$ , 623
- $f$ , 624
- $\imath$ , 367, 623
- $I$ , 632
- $j$ , 301, 623
- $\lambda$ , 485, 493, 495, 503, 505, 626
- $\pi$ , 178, 272, 462–465, 481, 486, 626
- $\psi$ , 240, 509, 626
- $\sigma$ , 485, 486, 626
- $\sqrt[3]{\cdot}$ , 622
- $\sqrt{\cdot}$ , 182, 622, 629
- $\times$ , 40, 128, 612, 627
- $V$ , 350, 522, 523
- $V_{\mathcal{N}}$ , 347, 523
- $-\mathcal{A}-$ , 33
- accuracy, *see* precision
- Achterberg, 284
- active set, 140, 244, 631
- adaptive, 301, 305, 307
- additivity, 84, 91, 495
- adjacency, 44
- adjoint  
    operator, 42, 43, 130, 447, 621, 626  
    self, 42, 351, 446, 485
- affine, 33, 182, *see also* function  
     $\emptyset$ , 33  
    algebra, 33, 34, 50, 65, 76, 579  
    boundary, 33, 39  
    combination, 56, 58, 64, 118  
    dimension, 21, 33, 50, 378–381, 387, 447, 627  
        complementary, 447  
        cone, 111  
        hyperplane, 63  
        low, 360  
        minimization, 471, 479  
        Précis, 381  
        reduction, 412  
        spectral projection, 480  
    face, 76  
    hull, *see* hull  
    independence, *see* independence  
    interior, 34  
    intersection, 70, 579  
    map, *see* transformation  
    membership relation, 172  
    normal cone, 619  
    set, 33, 50  
    subset, 33, 58  
        hyperplane intersection, 70  
        independence, 64, 70, 71  
        nullspace basis span, 70  
        parallel, 33  
        projector, 571, 579  
    transformation, *see* transformation
- a.i, 63, 64, 112, 623
- algebra  
    affine, 33, 34, 50, 65, 76, 579  
    arg, 173, 274, 529, 530  
    cone, 130  
    face, 76

- function
  - linear, 169
  - support, 184
- fundamental theorem, 503
  - linear, 69
- interior, 34, 39, 143
  - linear, 69, 485
  - nullspace, 515
  - optimization, 173, 529, 530
  - projection, 579
  - subspace, 32, 69, 515
- algebraic
  - complement, 77, 79, 122, 158, 574, 577, 624
  - projection on, 574, 600
  - multiplicity, 504
- Alizadeh, 238, 534
- alternating projection, *see* projection
- alternation, 405, 406, 481–483
- alternative, *see* theorem
- ambient, *see* vector, space
- anchor, 252, 253, 357, 363
- angle, 42, 60, 69, 634
  - acute, 60, 611
  - alternating projection, 611
  - brackets, 42, 43, 625
  - constraint, 350, 351, 353, 366
  - dihedral, 25, 69, 366, 410, 527, 634
  - EDM cone, 101
  - Gram definition, 348
  - halfspace, 60, 69
  - hyperplane, 69
  - inequality, 344, 409
    - matrix, 369
  - interpoint, 350
    - EDM definition, 348
  - obtuse, 60
  - positive semidefinite cone, 101
  - relative, 367
    - EDM definition, 369
  - right, 103, 124, 141, 142
  - torsion, 25, 366
- antidiagonal, 315, 372
- antisymmetry, 84, 116, 495, 632
- arg, 633
  - algebra, 173, 274, 529, 530
- artificial intelligence, 21
- asymmetry, 31
- audio, 28
  - compression, 330
  - echolocation, 20, 21
  - filter, 209–216
    - digital, 217
  - hearing, 602
    - compensation, 214, 215
  - metric, 215
  - op amp, 210, 211
  - sigma delta, 19, 20
  - signal dropout, 276–280
- Avis, 50
- axioms of the metric, 342
- axis of revolution, 102, 103, 141, 142, 527, 552
- $\mathcal{B}$  –, 174
- ball, *see* hypersphere
- 1-norm, 174, 175, 177
  - nonnegative, 280, 281
- 2-norm, 58, 604
  - Euclidean, 34, 174, 175, 354, 570, 603
    - smallest, 51
  - $\infty$ -norm, 173, 175, 177
  - nuclear norm, 55–57, 531
  - packing, 354
  - polyhedral, 175, 177
  - spectral norm, 58, 604
- Barker & Carlson, 114
- barrier, *see* interior-point
- Barvinok, 31, 90, 110, 112, 114, 227, 228, 238, 239, 242–244, 351, 353, 354, 363, 614
- base, 77
  - station, 364
- basis, 48, 149, 627
  - discrete cosine transform, 276
  - nonorthogonal, 149, 580
    - projection on, 580
  - nullspace, 61, 70, 73, 74, 100, 135, 151, 160, 507, 575, 583, 584, 628
    - Schur-form, 500
    - vectorized, 74, 100
  - orthogonal, 48, 49, 73, 577
    - projection on, 580
  - overcomplete, 74, 627
  - pursuit, 299
  - range, 87, 507, 575, 580, 627
    - complement, 151
    - EDM, 422
    - vectorized, 74
  - rowspace, 627
  - standard
    - matrix, 48, 49, 590, 614, 622
    - vector, 48, 120, 349, 632
- Bauschke, 610
- bees, 27, 354
- Bellman & Fan, 223
- Ben-Israel, 134, 230
- Bessel, 576
- bijection, 631, *see* bijective
- bijective, 42, 44, 46, 94, 145, 631
  - Gram-form, 375, 376
  - inner-product form, 377, 378
  - invertible, 44, 568
  - isometry, 45, 406
  - linear, 43, 45, 47, 49, 112
  - orthogonal, 45, 406, 465, 525, 538
- binary
  - distance matrix, 202
  - program, *see* program
  - search, 206, 212–214
- binomial, *see* coefficient
- biorthogonal
  - decomposition, 571, 577, 578
  - expansion, *see* expansion
- biorthogonality, 147, 149
  - condition, 145, 504, 518
- Birkhoff, 57
- bisection method, 206, 212–214
- Biswas
  - Pratik, 254, 266, 366
  - Suddhendu, 515
- blades, 127

- Blu, 299  
 Blumenthal, 342  
 Boone, 585  
 Borg & Groenen, 25, 404  
 Borwein, 74, 84, 595, 610  
 bound  
     greatest  
         lower, 231, 541, 633  
         upper, 226, 351, 463  
     least  
         lower, 275, 360, 458, 463, 464  
         upper, 633  
     polyhedral, 418  
 boundary, 34, 35, 39, 53, 76, 629  
 0, 82, 85  
 extreme, 88  
 of affine, 33, 39  
 of cone, 85, 143–145  
     EDM, 416, 418, 422, 428  
     membership, 138, 231  
     membership relation, 132, 138, 231  
     over PSD, 143  
     polyhedral, 145  
     PSD, 35, 91, 96, 107–109, 231, 432  
     ray, 85  
 of halfspace, 60  
 of hypercube slice, 307  
 of open set, 35  
 of point, 39  
 of ray, 77  
 -point, 78  
 relative, 34, 35, 39, 53, 76  
 bounded, 51  
 bowl, 192, 198, 551, 552  
 box, 67, 603, *see also* hypercube  
 Boyd, 469, 473  
 bridge, 366, 387, 439, 458, 470  
 Bronstein, 24  
 Brooks, 585  
 Bunt, 109, 595
- C –, 31
- calculus  
     matrix, 543  
     of inequalities, 7  
     of variations, 135  
 Carathéodory's theorem, 131, 590  
 card, 634  
 cardinality, 271, 346, 634, 635  
     1-norm heuristic, 272, 273, 472  
     constraint, *see* constraint  
     convex envelope, *see* convex  
     geometry, 175, 177, 282  
     minimization, *see* minimization  
     monotonicity, 196  
     -one, 175, 177, 281, 282  
     quasiconcavity, 179, 196, 271  
     reduction, *see* minimization  
     regularization, 277, 298  
 Carlson, 114  
 Cartesian  
     axes, 33, 52, 65, 80, 86, 121, 126, 273, 603  
     cone, *see* orthant  
     coordinates, 403, 631
- product, 40, 228, 612, 627  
     cone, 81, 82, 128  
     function, 172  
     subspace, 273, 603  
 cartography, 26, 402–404  
     stellar, 19, 367, 407  
 Cauchy-Schwarz, 71, 610  
 Cayley-Menger  
     determinant, 410, 412  
     form, 382, 396, 463, 480  
 cellular, 365  
     telephone network, 364  
 center  
     geometric, 349, 352, 370, 401, 627  
     gravity, 370  
     mass, 370  
 central path, 228  
 certificate, 168, 230, 231, 248, 315, 336  
     null, 132, 139  
 Chabrillac & Crouzeix, 434  
 chain rule, 548  
     two argument, 548  
 Chu, 7  
 c.i, 111, 112, 114, 623  
 circle  
     antenna, 22  
     conic section, 102–104  
     Fantope, 54  
     Fermat point, 352  
     hyperplane, 62  
     intersection, 289, 357, 358  
     optimization, 529  
     pyramid, 412  
     quarter, 76  
     singular value decomposition, 507, 508  
     unit, 217  
 clipping, 180, 456, 603, 604, 625  
 closed, *see* set  
 coefficient  
     binomial, 178, 284, 625  
     projection, 46, 431, 433, 498, 575, 588–591  
         nonorthogonal, 147, 587, 588  
 cofactor, 411  
 combinatorial, *see* problem  
 compaction, 254, 266, 350, 353, 436, 473  
 comparable, 84, 146  
 complement  
     algebraic, 77, 79, 122, 158, 574, 577, 624  
         projection on, 574, 600  
         projector, 574  
     halfspace, 58  
     orthogonal, 46, 69, 250, 574, 577, 624, 628  
         cone, dual, 122, 125, 126, 600  
         cone, normal, 619  
         projection on, 600  
         vector, 628  
     relative, 74, 624  
 complementarity, 158, 173, 272  
     linear, 161  
     maximal, 223, 228, 243  
     problem, 161  
         linear, 161  
     semidefinite, 223, 233, 248  
     strict, 233  
 complementary

- affine dimension, 447  
 dimension, 69  
 eigenvalues, 574  
 halfspace, 61  
 inertia, 397, 499, 500  
 slackness, 233  
 subspace, *see* complement  
 complexity, *see* computational intensity  
 compressed sensing, 173–177, 179, 238, 274, 275, 285, 286, 299–301, 306, 330  
 formula, 275, 299  
 nonnegative, 176, 177, 275, 280–284  
 Procrustes, 290  
 prototypical, 280  
 compression, signal, 330  
 compressive sampling, *see* compressed  
 computational intensity  
 alternating projection, 617  
 $O$ , 632  
 presolver, 284  
 rank reduction, 239  
 SDP, 223, 254, 266  
 concave, 169, 170, 180  
 condition  
 biorthogonality, 145, 504, 518  
 convexity, 197  
 first-order, 197, 199, 201, 203  
 second-order, 200, 204, 205  
 Moore-Penrose, 567, 569  
 necessary, 624  
 optimality, *see* optimality  
 orthogonality, 575  
 orthonormality, 148, 149, 576  
 sufficient, 624  
 cone, 58, 77, 81, 629, 634  
 $\emptyset$ , 77  
 $0$ , 82, 85, *see* cone, origin  
 algebra, 130  
 blade, 127  
 boundary, *see* boundary  
 Cartesian, *see* orthant  
 product, 81, 82, 128  
 circular, 37, 81, 102–105, 181, 451  
 axis, 102–104  
 Lorentz, *see* cone, Lorentz  
 right, 102  
 closed, *see* set, closed  
 convex, 58, 77, 81, 82, 87  
 subsets, 108  
 difference, 40, 130  
 dimension, 111  
 dual, 122–128, 136, 151–157, 162–164, 629  
 Cartesian product, 128  
 construction, 123, 124  
 convex, 122  
 dual, 129, 136  
 facet, 145, 165  
 formula, 136, 151–157, 163  
 full-dimensional, 129  
 halfspace-description, 122, 123, 136, 137  
 in subspace, 130, 132, 149, 152, 153  
 intersection, 130  
 of intersection, 130  
 of nonsimplicial, 162  
 origin, 125  
 orthogonal complement, 122, 125, 126, 600  
 pointed, 119, 129  
 product, Cartesian, 128  
 proper, 130  
 property, 128  
 self-, 132, 141, 142  
 spectral, 398, 399  
 unique, 122, 145  
 vertex-description, 151, 162–165  
 $\sqrt{\text{EDM}}$ , 417, 419, 468  
 EDM, 346, 415–417, 435  
 angle, 101  
 axis, 101, 394, 417, 423, 424  
 boundary, 416, 418, 422, 428  
 circular, 102  
 construction, 423  
 convexity, 419  
 decomposition, 448  
 dual, 441, 442, 446, 449–452  
 equality, 26, 415, 443, 450  
 extreme direction, 425  
 face, 424, 426  
 interior, 416, 422  
 one-dimensional, 88, 418, 432, 441  
 origin, 418, 441  
 polar, 430, 448, 449, 479, 480  
 positive semidefinite, 426, 432, 443  
 face, 85–87  
 intersection, 138  
 pointed, 82  
 PSD, 116  
 smallest, 82, 138–140, 177, 284  
 smallest, EDM, 424  
 smallest, generators, 139, 140, 284  
 smallest, PSD, 97–100, 110, 431, 443, 450  
 smallest, subspace, 98  
 facet, 145, 165  
 full-dimensional, 77, 86, 119, 125, 128, 140, 141, 145, 150, 155  
 dual, 129  
 halfline, 77, 81, 82, 87, 88, 112, 119  
 halfspace, *see* halfspace  
 hull, 88, 101, 130  
 ice cream, *see* cone, circular  
 interior, *see* interior  
 intersection, 81  
 dual of, 130  
 of dual, 130  
 pointed, 82  
 invariance, *see* invariance  
 line, 82, 119, 130  
 Lorentz  
 circular, 81, 103, 181  
 dual, 125, 141  
 face, 116  
 polyhedral, 117, 125  
 majorization, 488  
 membership, *see* membership relation  
 monotone, 155–157, 397, 482, 539, 604  
 nonnegative, 153–155, 398, 399, 403, 405, 406, 463, 464, 482, 488, 604  
 negative, 146, 581  
 nonconvex, 77–80  
 nonpointed, 82, 87, 145  
 c.i., 111

- dual, 130
- EDM, dual, 441
- monotone, 155, 156, 397
- polyhedral, 86, 114, 115, 119, 128
- spectral, 398, 399
- nonsimplicial, 112, 113, 162, 164, 165
- normal, 158–161, 172, 428, 596, 601, 607, 617–620, 628
- affine, 619
- elliptope, 618
- membership relation, 158
- origin, 158
- orthant, 158
- orthogonal complement, 619
- translated, 607, 617, 619, 620
- one-dimensional, 77, 81, 82, 85, 112, 117, 119
- EDM, 88, 418, 432, 441
- PSD, 87
- origin, 58, 77, 81, 82, 85, 112, 117
- dual, 125
- EDM, 418, 441
- normal, 158
- pointed, 87, 119
- orthant, *see* orthant
- pointed, 82–87, 119, 131, 144, 147, 157
- closed convex, 82–85, 114, 117
- dual, 119, 129
- face, 82
- intersection, 82
- nonconvex, 82
- nullspace, 119
- origin, 87, 119
- polyhedral, 86, 119, 128, 144
- polar, 77, 122, 428, 479, 598–601, 604, 629
- EDM, 430, 448, 449, 479, 480
- positive semidefinite, 445, 449
- polyhedral, 58, 112–120
- dual, 113, 130, 132, 136, 142, 144, 145
- equilateral, 142
- face, 115, *see* cone
- facet, 145, 165
- halfspace-description, 116
- majorization, 482
- nonpointed, 86, 114, 115, 119, 128
- pointed, 86, 119, 128, 144
- proper, 113
- vertex-description, 58, 119, 132, 144
- positive semidefinite, 37, 90–92, 445
- angle, 101
- axis, 102
- boundary, 35, 91, 96, 107–109, 231, 432
- circular, 102–105, 181
- convex subsets, 108
- dimension, 97
- dual, 141, 153
- EDM, 426, 432, 443
- equality, 26, 415, 443
- extreme direction, 87, 92, 101
- face, *see* cone, face
- hyperplane, supporting, 97, 98
- inscription, 106
- interior, 90, 230
- inverse image, 93, 144, 187, 204, 209, 228
- one-dimensional, 87
- optimization, 227
- polar, 445, 449
- polyhedral, 105
- rank, 97–100, 108, 224, 496
- visualization, 225
- product, Cartesian, 81, 82, 128
- proper, 85
- dual, 130
- nonsimplicial, 162
- polyhedral, 113
- quadratic, 81, 181
- ray, 77, 81, 82, 112, 119
- EDM, 88
- positive semidefinite, 87
- recession, 125
- rotated, 181
- EDM, 451
- orthant, 120, 142, 148
- positive semidefinite, 105
- second-order, *see* cone, Lorentz
- selfdual, *see* cone, dual
- shrouded, 142, 452
- simplicial, 59, 112, 120, 121, 146, 604
- decomposition, 162–164, 412
- dual, 130, 146, 152
- spectral, 396–398, 462, 626, 629
- dual, 398, 399
- orthant, 464
- subspace, 81, 117, 119
- sum, 81, 130
- supporting hyperplane, *see* hyperplane
- tangent, 428
- transformation, *see* transformation
- two-dimensional, 112
- union, 77, 79, 130, 163
- unique, 90, 122, 416
- vertex, 82
- none, 82, 86, 87, 114, 115, 119, 128
- congruence, 373, 623
- transformation, *see* transformation
- conic
- combination, 58, 77, 118
- constraint, 160, 162, 180, 181
- coordinates, 122, 149, 151, 165–168
- inversion, 166
- hull, *see* hull
- independence, *see* independence
- problem, 127, 160, 162, 172, 222
- dual, 127, 222
- section, 102
- circular cone, 104
- conjugate
- complex, 302, 488, 512, 525, 591, 621, 625
- convex, 129, 131, 132, 136, 145, 152, 158
- function, 472, 625
- gradient, 305
- conservation
- dimension, 69, 380, 510, 511, 521
- constellation, 19, 21, 367
- constraint
- angle, 350, 351, 353, 366
- binary, 234, 292, 326
- cardinality, 271–275, 285, 286, 304
- full-rank, 238
- nonnegative, 176, 179, 271–275, 280–284
- projection on PSD cone, 297

- rank-one, 297
- conic, 160, 162, 180, 181
- equality, 135, 162
  - affine, 221
- Gram, 350, 351, 353, 366
- inequality, 161, 244
- log, 191
- nonnegative, 57, 71, 72, 161, 173–177, 222, 286
- norm, 286, 289, 291, 293, 510, 585
  - Frobenius, 586
  - spectral, 58
- orthogonal, 335, 367, *see* Procrustes
- orthonormal, 287
- permutation, *see* polyhedron, permutation
- polynomial, 180, 181, 291
- quadratic, 224, 291
  - convex, 190
  - nonconvex, 190
- qualification, *see* Slater
- rank, 247–251, 287, 291, 332, 333, 474, 480
  - cardinality, 297
  - indefinite, 329
  - projection on matrices, 604
  - projection on PSD cone, 193, 297, 462, 466, 532
  - wide, 329
- Schur-form, 58
- singular value, 58
- sort, 405
- tractable, 221
- content, 411, 634
- contour plot, 159, 199
- contraction, 58, 161, 305, 307
- control theory, 183, 483, 616
- convergence, 250, 406, 606
  - geometric, 607
  - measure, 236
- convex, 7, 19, 31, 169, 171, 631
  - combination, 31, 58, 118
  - envelope, 471, 472
    - bilinear, 209
    - cardinality, 237, 472
    - rank, 261, 363, 471, 472
  - form, 8, 221
  - geometry, 31
    - fundamental, 60, 65, 68, 341, 349
  - hull, *see* hull
  - iteration
    - accelerant, 338
    - cardinality, 271–276, 285
    - cardinality-one, 281, 282
    - convergence, 250
    - indefinite, 329
    - optimality, 248, 250, 251, 272
    - optimality, local, 251, 273, 276, 285, 336
    - rank, 247–251, 474
    - rank-one, 287, 333–336, 483
    - stall, 251, 273, 276, 287, 293, 333
    - wide, 329
  - log, 218
  - optimization, 7, 19, 171, 221, 357
    - art, 8, 238, 254, 484
    - fundamental, 221
    - solution set, 171, 333
  - polyhedron, *see* polyhedron
  - problem, *see* problem
  - program, *see* problem
  - quadratic, *see* function
  - strictly, *see* function, convex
  - convexity, 169, 180
    - condition
      - first-order, 197, 199, 201, 203
      - second-order, 200, 204, 205
    - eigenvalue, 387, 497, 532
    - K-, 170
    - log, 218
    - norm, 58, 170–175, 191, 196, 530, 635
      - of difference, 346
    - projector, 574, 577, 595, 596
    - property, 218
    - strict, *see* function, convex
  - coordinates, conic, *see* conic
  - cosine, *see* trigonometry
  - CPU, 630
  - Crippen & Havel, 366
  - Critchley, 376, 479
  - criterion
    - EDM, 349, 350, 430, 434
      - dual, 446
    - Finsler, 434
    - metric *versus* matrix, 382
    - Schoenberg, 25, 26, 349, 350, 396, 416, 427, 430, 447, 450–453, 591
  - Crouzeix, 434
  - cube, 544, 546, 631
  - curvature, 200, 207
  - cvx, 265, 291, 327
  - cylinder, 55
  - $D$  –, 349
  - Dahl, 328
  - Dantzig, 57, 222, 223
  - d'Aspremont, 92, 297
  - dB, 630
  - DC, 630
  - DCT, *see* discrete cosine transform
  - decomposition, 631
    - biorthogonal, 571, 577, 578
    - completely PSD, 81, 251
    - dyad, 514
    - eigenvalue, 503–505, 509
    - extreme direction, 102
    - factorization
      - nonnegative, 251
      - positive, 251
      - spectral, 266–268
    - Moreau, 161, 600
    - nullspace, 61, 570
    - orthonormal, 576, 577
    - simplicial, 162–164, 412
    - singular value, *see* singular value
  - deconvolution, 299
  - definite, *see* matrix, positive
  - Definitions, 621
  - de Leeuw, 399, 459, 469, 475
  - Delsarte, 354
  - deprecation, 140, 282
  - dereverberation, 20, 21
  - derivative, 193, 629

- directional, 200, 549, 551, 552, 555, 556
  - dimension, 551, 558
  - gradient, 552
  - second, 553
- gradient, 556
  - multivariate, 553
  - partial, 543, 549, 629
  - table, 558–565
  - trace, 205, 562
- D'Errico, 503
- description
  - conversion, 122
  - halfspace-, 59–61
    - of affine, 70
    - of dual cone, 122, 123, 136, 137
    - of polyhedral cone, 116
    - of polyhedron, 116
    - of subspace, 69
  - vertex-, 58, 59, 87
    - of affine, 70
    - of dual cone, 151, 162–165
    - of halfspace, 114, 115
    - of hyperplane, 63
    - of polyhedral cone, 58, 119, 132, 144
    - of polyhedron, 53, 118
    - of subspace, 69, 118
    - projection on affine, 584
- determinant, 206, 408, 493, 502, 558, 564
  - Cayley-Menger, 410, 412
  - inequality, 496
  - product, 494
- Deutsch, 596, 603, 604, 607
- Deza, 240, 430
- DFT, *see* discrete Fourier transform
- diag(), *see*  $\delta$
- diagonal, 485–487, 495, 569, 626
  - $\delta$ , 485–487, 493, 495, 503, 569, 626
  - binary, 572
  - commutative, 494
  - diagonalization, 505
  - dominance, 107
  - inequality, 496
  - nonnegative, 495
  - positive semidefinite, 495
  - pseudoinverse, 569
  - values
    - eigen, 503
    - singular, 506
    - zero, 510, 569
- diagonalizable, 148, 503–505
  - simultaneously, 98, 233, 493, 494, 513
- diagonalization, 37, 74, 503–505
  - diagonal matrix, 505
  - expansion by, 148
  - symmetric, 505
- diamond, 175, 372
- difference, 630
  - cone, 40, 130
  - of functions, 179, 190, 274
  - positive semidefinite, 110
  - set, 40, 130, 624
  - vector, 40, 130, 624
- differentiable, *see* function
- differential, 629
  - discrete, 630
- Fréchet, 550
- Gâteaux, 550
  - partial, 629
- diffusion, 436
- dilation, 404
- dimension, 33, 634
  - affine, *see* affine
  - complementary, 69
  - conservation, 69, 380, 510, 511, 521
  - domain, 45
  - embedding, 50
  - Euclidean, 627
  - face, 75
    - positive semidefinite cone, 97, 98
  - invariance, 45
  - nullspace, 69
  - Précis, 381
  - range, 45
  - rank, 111, 380, 381, 634
- Dirac, 299, 515
- direction, 77, 86, 549
  - extreme, *see* extreme
  - matrix, 248–251, 261, 298, 355, 475
    - analytical, 249
    - existence, 248
    - Identity, 250, 330
    - interpretation, 248
  - Newton, 193
  - projection, 575, 577, 587, 588
    - nonorthogonal, 572, 573, 582, 587
    - parallel, 578, 579, 581, 582
  - steepest descent, 193, 551
  - vector, 248–251, 271–274
    - analytical, 273
    - existence, 248
    - Identity, 250, 330
    - interpretation, 272, 273, 472
    - optimality condition, 248, 272
    - unity, 272
- disc, 53
- discrete
  - cosine transform, 276, 301
  - Fourier transform, 45, 301, 306, 622, 630
    - inverse, 301, 303
    - time Fourier transform, 217
  - discretization, 136, 170, 200, 450
  - dissimilarity, 400
  - dist, 595, 608, 634
  - distance, 346, 608, 609
    - 1-norm, 174, 175, 202, 203, 394
    - absolute, 342, 394
    - binary, 202
    - comparative, *see* distance, ordinal
    - Euclidean, 346, 367, 368
    - geometry, 19, 364
    - matrix, *see* EDM
    - maximization, 597
    - minimization, 575, 592
    - ordinal, 402–407
    - origin to affine, 63, 570, 583
      - 1-norm, 174, 175
      - nonnegative, 176, 177, 281
    - property, 342
    - self, 342
    - square, 346, 348

- taxicab, *see* distance, 1-norm  
 distortion, 210, 300, 330  
     spatial, 47  
 Dokmanić, 20  
 Donoho, 330  
 doublet, *see* matrix  
 dropout problem, 276–280  
 dual  
     affine dimension, 447  
     dual, 129, 132, 136, 232, 235, 296  
     feasible set, 228  
     function, 127, 129, 231, 296  
     norm, 125, 530, 635  
     problem, 125–128, 167, 178, 222, 231, 234, 441  
         strong, 128, 129, 232, 296, 351, 538  
         via primal, 125, 127–129, 234  
     projection  
         on cone, 77, 122, 449, 479, 598–601  
         on convex set, 595, 597, 598  
         on EDM cone, 449, 479, 480  
         on subspace, 599, 600  
     variable, 125, 161, 162, 167, 625  
 duality, 125, 351  
     gap, 128, 129, 229, 232, 233, 296  
     strong, 127–129, 167, 183, 233  
     weak, 127, 231  
 Duensing, 301  
 Dulá, 585  
 dvec, 49, 405, 634  
 D:Wave, 323–327  
 dyad, *see* matrix  
 Dykstra, 608, 617  
     algorithm, 606, 616, 617, 619, 620  
  
     –  $\mathcal{E}$  –, 388  
 earO, 214, 215  
 Ebert, 567  
 Eckart & Young, 27, 460, 462, 484, 532  
 edge, 75  
 $\sqrt{\text{EDM}}$ , 419  
 EDM, 26, 341, 346, 453, 630  
     closest, 461  
     composition, 392–394, 417–419  
     construction, 421  
     criterion, 349, 350, 430, 434  
         dual, 446  
     definition, 346, 419  
         Gram-form, 348, 428  
         inner-product form, 367  
         interpoint angle, 348  
         relative-angle form, 369  
     dual, 446, 447  
 eigenvalue, 381, 395, 404, 420  
 exponential entry, 392  
 graph, 344, 354, 358, 360, 361  
 indefinite, 395  
 invariance, *see* invariance  
 membership relation, 450  
 nonnegative, 446  
 projection, 433  
 range, 422  
 rank, 381, 422  
 subgraph, 436, 437  
 test, numerical, 344, 349, 460  
  
     unique, 391, 401, 462  
 eigen, 503  
     matrix, 91, 148, 504  
     spectrum, 396  
         ordered, 397, 462–464  
         unordered, 398  
     value, 491, 503, 530, 532  
          $\lambda$ , 485, 493, 495, 503, 505, 626  
         convexity, 387, 497, 532  
         decomposition, 503–505, 509  
         distinct, 503, 504  
         EDM, 381, 395, 404, 420  
         Identity, 492, 532  
         inequality, 493, 537  
         interlaced, 386, 395, 497, 528  
         intertwined, *see* eigenvalue, interlaced  
         inverse, 504, 505, 510  
         largest, 387, 497, 532–534  
         left, 503  
         maximum, *see* eigenvalue, largest  
         minimum, *see* eigenvalue, smallest  
         of sum, 497  
         positive semidefinite, 492, 496, 509  
         precision, 503  
         principal, 297, 532  
         product, 493–495, 509, 535  
         projection, 572, 574, 576, 588–590  
         real, 488, 503, 505  
         repeated, 109, 503, 504  
         Schur, 501  
         singular value, 506, 509, 510  
         smallest, 387, 393, 441, 497, 532, 533  
         sum of, 44, 485, 493, 533, 534  
         symmetric matrix, 505, 509  
         transpose, 503  
         unique, 102, 503  
         zero, 510, 516  
     vector, 503, 532, 533  
         distinct, 504  
         EDM, 420  
         left, 503, 504  
         Li, 503, 504  
         normal matrix, 505  
         orthogonal, 504, 505  
         principal, 297, 298, 532, 586  
         real, 503, 505, 511  
         symmetric matrix, 505  
         unique, 504  
     elbow, 457, 458  
     element, 631  
     El Ghaoui, 297  
     ellipse, singular value decomposition, 507, 508  
     ellipsoid, 31, 36, 37, 286  
         invariance, 42  
     elliptope, 235, 237, 350, 357, 388, 389, 394, 428, 429  
         smooth, 388  
         vertex, 236, 388, 389  
     embedding, 378  
         dimension, 50  
     empty  
         interior, 34  
         set, *see* set  
     entanglement, 324  
     entry, 33, 170, 621, 625, 626, 629, 631  
         largest, 178

- smallest, 178
- zero, 510
- epigraph, 169, 171, 185, 204, 218
  - form, 189, 190, 478, 586
  - intersection, 185, 191, 219
  - nonconvex, 185
- equal loudness, 214, 215
- equality
  - constraint, 135, 162
    - affine, 221
  - EDM cone, 26, 415, 443
    - dual, 450
  - PSD cone, 26, 415, 443
- equation, normal, 569
- equivalent, 631, *see* problem
- Ericson, 356
- errata, 164, 183, 252, 376, 534, 602
- Eternity II, 307–321
  - extreme point, 316, 317
  - hypersphere, 321
  - maximization, 321
  - quantum, 326, 327
- Euclidean
  - ball, 34, 174, 175, 354, 570, 603
    - smallest, 51
  - distance, 341, 346
    - geometry, 19, 358, 364, 591
    - matrix, *see* EDM
    - metric property, 342, 383
      - fifth, 343–345, 407–409, 413
    - norm, *see* 2-norm
    - projection, 108
    - space, 31, 32, 341, 627
      - ambient, 32, 441, 451
- exclusive
  - mutually, 134
- expansion, 149, 631
  - biorthogonal, 122, 131, 145, 148, 149, 572, 578
    - EDM, 425
    - projection, 577, 580–582
    - unique, 147–150, 426, 518, 577, 580–582
  - implied by diagonalization, 148
  - orthogonal, 48, 122, 148, 576, 577
  - w.r.t orthant, 149
- exponent, *see* function, fractional
- exponential, *see* matrix
- exposed, 74, 75
  - closed, 426
  - direction, 88
  - extreme, 78, 82, 115
  - face, 75, 116
  - point, 75, 78, 88
    - density, 76
- extreme, 74
  - boundary, 88
  - direction, 72, 86–89
    - conic independence, 114
    - distinct, 86
    - dual cone, 113, 124, 145, 149–155, 163, 164
    - dual monotone cone, 155
    - EDM cone, 88, 425, 446
    - monotone nonnegative cone, 154, 155
    - none, 82, 86, 87, 114, 115, 119, 128
    - of  $\mathbf{0}$ , 87
    - one-dimensional, 87, 88
- orthant, 72
- PSD cone, 87, 92, 101
- shared, 165
- supporting hyperplane, 145
- unique, 86
- exposed, 78, 82, 115
- point, 57, 74–76, 78, 227
  - cone, 82, 110
    - Fantope, 249
    - feasible set, 248
    - hypercube slice, 273
  - ray, 86, 428
- sum, 87
- transitivity, 76
- $-\mathcal{F}-$ , 76
- $f$ , 624, *see* function
- face, 36, 75
  - $\emptyset$ , 75
  - affine, 76
  - algebra, 76
  - cone, *see* cone
  - exposed, *see* exposed
  - halfspace, 75
  - hyperplane, 75
  - intersection, 138
  - isomorphic, 97, 351, 424
  - of face, 76
  - polyhedron, 115, 116
  - smallest, 76, 97–100, 138–140, 622, *see* cone
  - subspace, 75
  - transitivity, 76
- facet, 75, 145, 165
- facial recognition, 24
- factorial, 625
- factorization, *see* decomposition
- Fan, 53, 223, 248, 261, 493, 530, 533, 537
- Fantope, 53, 54, 102, 188, 248, 249
  - extreme point, 249
  - Linear Program analogue, 272
- Farkas' lemma, 131–134
  - positive definite, 230
  - not, 231
  - positive semidefinite, 229
- fax, 203
- Fazel, 473
- feasibility, *see* problem
- feasible, *see* solution
- Fejér, 141, 196, 610, 613
- Fenchel, 23, 169, 472
- Fermat point, 352
- Ferreira, 604
- Feynman, 324
- Fiacco & McCormick, 223
- fifth metric property, *see* Euclidean
- filter
  - bank, 602
  - design, 484
    - arbitrary magnitude, 209, 215, 217
    - implementation, 267
- find, 631, 633
- finitely generated, 116, 117
- Finsler, 434
- floor, 635

- flowgraph, 211  
 Forsgren, Gill, & Wright, 221  
 Fourier transfer function, 266  
 Fourier transform, 266, 267  
     discrete, 45, 301, 306, 622, 630  
     inverse, 301, 303  
     -time, 217  
     fast, *see* discrete  
 frame, 114  
 Frankel, Felice, 691  
 Fréchet differential, 550  
 Frisch, 223  
 Frobenius, 43, *see* norm  
 Fukuda, 50  
 full, 34  
     -dimensional, 34, 85  
     cone, 77, 86, 119, 125, 128, 140, 141, 145, 150, 155  
     cone, dual, 129  
     -rank, 70, 238  
 function, 169, 624, *see operator and transformation*  
     affine, 40, 42, 172, 182–184, 194, 195  
         inverse, 40, 42  
         monotonic, 196  
         supremum, 185, 530  
         transformation, 196, 218  
     bilinear, 209  
     biquadratic, 209  
     composition, 180, 196, 219, 548  
         affine, 196, 218  
     concave, 169, 170, 180, 473, 532  
     conjugate, 472, 625  
     continuity, 170, 185, 201, 218  
     convex, 169, 185, 347, 369  
         difference, 179, 190, 274  
         invariance, 170, 185, 218  
         nonlinear, 170, 184  
         simultaneously, 187, 188, 204  
         strictly, 170, 171, 196–198, 200, 205, 469, 473, 477, 551, 552  
         sum, 171, 196, 218, 219  
         supremum, 185, 191, 219  
     differentiable, 170, 184, 197, 199, 205  
     non, 170, 174, 194, 218, 305, 306  
     distance  
         1-norm, 202, 203  
         binary, 202  
         Euclidean, 346, 367, 368  
     dual, *see* dual  
     fractional, 187, 197, 201, 204  
         exponent, 181, 565  
         inverted, 180, 181, 558, 565  
         maximization, 469  
         minimization, 500  
         power, 181, 565  
         projector, 188  
         pseudo, 186  
         root, 180, 622  
         square root, 182, 469, 496, 505, 622  
     inverted, *see* function, fractional  
     invertible, 42, 44, 45  
     Lagrangian, *see* Lagrangian  
     linear, *see* operator  
     log, *see* log  
     matrix, 169, 171, 201  
     chain rule, 548  
     convexity, 203, 205  
     epigraph, 204  
     gradient, 545  
     line theorem, 205  
     product, 546  
     monotonic, 169, 178, 195, 216  
         affine, 196  
         bilinear, 209  
         biquadratic, 209  
         cardinality, 196  
         composition, 196, 219  
         norm, 173, 179, 274  
         quasilinear, 208, 218  
         rank, 196, 471  
         strictly, 195, 208  
         multidimensional, 169, 201, 632  
         affine, 172, 183  
         convexity, 203, 204  
         gradient, 192  
         line theorem, 205  
         monotonic, 195  
         quasiconvex, 207  
         negative, 180  
         nonconvex, 185, 201, 458, 475  
         norm, 173, *see* norm  
         objective, *see* objective  
         odd, 240, 509  
         presorting, 178, 272, 462–465, 481, 486, 626  
         product, 187, 197, 208, 209, 216  
         projection, *see* projector  
         quadratic, 200, 205, 487, 500, 591  
         binary, 326, 327  
         convex, 171, 190, 347, 368, 370, 540  
         convex, strictly, 171, 192, 198, 200, 205, 477, 551, 552  
         distance, 346  
         maximization, *see* maximization  
         minimization, *see* minimization  
         nonconvex, 291, 540  
         nonnegative, 500  
         quasiconcave, 207–209, 216, 218, 240  
         cardinality, 179, 196, 271  
         not, 218  
         rank, 107  
         strictly, 207  
         quasiconvex, 108, 185, 206–209, 216, 218, 368  
         continuity, 207, 218  
         not, 218  
         strictly, 207  
         quasilinear, 201, 208, 218, 509  
         quotient, *see* function, fractional  
         ratio, *see* function, fractional  
         signum, 208, 240, 305, 509, 558, 559, 565, 626, 633  
         smooth, 170, 174, 306  
         sorting, 178, 272, 462–465, 481, 486, 626  
         square, *see* norm  
         square root, *see* function, fractional  
         sstress, 458, 459  
         step, 626  
             matrix, 240  
             vector, 509  
         strain, 459  
         stress, 201, 458, 459, 475

- support, 67, 184, 596  
     algebra, 184  
 transfer, 209–211, 217, 266–268  
 trivial, 171  
 vector, 169, 171, 626  
     convexity, 199  
     epigraph, 185  
     gradient, 543, 545  
     line theorem, 205  
     monotonicity, 196  
     sublevel, 186  
 fundamental  
     convex  
         geometry, 60, 65, 68, 341, 349  
         optimization, 221  
     metric property, 342, 383  
     semidefiniteness test, 141, 488, 591  
     subspace, 41, 69, 516, 519, 520, 568, 575  
         projector, 575
- $\mathcal{G}$  –, 137
- Gaffke & Mathar, 442  
 Gale, 360  
 Galtier, 181  
 gap  
     bisection, 212  
     duality, 128, 129, 229, 232, 233, 296  
     quantum, 326–328  
     time, 276–280  
     trace–rank, 471  
 Gâteaux differential, 550  
 generators, 50, 58, 87, 118, 137, 622  
     c.i., 114  
     face, smallest, 139, 140, 284  
     finite, 116, 117  
     hull, 118  
         affine, 50, 63, 584  
         conic, 58, 114, 115  
         convex, 53  
     l.i., 114  
     minimal set, 53, 63, 114, 320, 627  
         affine, 584  
         extreme, 87, 145  
         halfspace, 114, 115  
         hyperplane, 114  
         orthant, 170  
         PSD cone, 92, 141, 202  
     unique, *see* unique
- geometric  
     center, 349, 352, 370, 401, 627  
     operator, 376  
     subspace, *see* subspace
- Hahn-Banach theorem, 60, 65, 68  
 mean, 565  
 multiplicity, 504
- Geršgorin, 105  
 Gill, 221  
 gimbal, 526  
 global  
     optimality, *see* optimality  
     positioning system, 21, 253
- Glossary, 621  
 Glunt, 477, 617  
 Golub, 489
- Gordan, 134  
 Gould, 223  
 Gower, 341, 387  
 GPS, *see* global  
 gradient, 135, 183, 192, 193, 197, 198, 543, 551, 557  
     affine, 184, 194, 195  
     composition, 548  
     conjugate, 305  
     derivative, 556  
         directional, 552  
     dimension, 192  
     first-order, 556  
     image, 304  
     monotonic, 196  
     norm, 179, 180, 192, 559, 561, 575, 583, 584, 592, 595  
     normal, 159, 193, 198  
     of nonlinear  $f$ , 184  
     product, 546  
     second-order, 557  
     sparsity, 300  
     table, 558–565  
     zero, 135, 193, 194, 592
- Gram  
     -form  
         bijective, 375, 376  
         EDM definition, 348, 428  
         injective, 371, 375, 376  
         invariance, 371  
         inversion, 376  
         problem, *see* problem, proximity  
     matrix, *see* matrix
- Gregory, 354  
 Groenen, 25, 404  
 Gross, 27
- $\mathcal{H}$  –, 60
- Hadamard  
     product, 43, 468, 485–487, 548, 623  
         of vectors, 548  
         positive semidefinite, 495  
         transpose, 487  
         quotient, 558, 565, 623  
         square root, 622
- Hahn-Banach, 60, 65, 68  
 halfline, 77, 81, 82, 87, 88, 112, 119, 209  
 halfplane, 125, 127, 209, 267  
 halfspace, 58–61, 114, 115, 119, 128  
      $\mathcal{H}_+$ , 60, 629  
      $\mathcal{H}_-$ , 60, 629  
     angle, 60, 69  
     boundary, 60  
     complement, 58  
     -description, *see* description  
     interior, 39  
     intersection, 61, 129, 141  
         cone, 116  
         dual cone, 122  
         polyhedra, 116  
     polarity, 60  
     vertex-description of, *see* description
- halftone, 203  
 Han, 616  
 Hardy-Littlewood-Pólya, 463

- Havel, 366  
 Hayden & Wells, 189, 415, 420, 426, 476, 477  
 hearing  
     aid, 602  
     compensation, 214, 215  
     test, 214  
 Herbst, 509  
 Hessian, 193, 543  
 hexagon, 353  
 Hindi, 473  
 Hiriart-Urruty, 67, 541  
 hollow, 622, *see* matrix & subspace  
 homogeneity  
     convexity, 218  
     EDM, 347  
     norm, 173  
     partial order, 84  
     projection, 600  
 homotopy, 272  
 honeycomb, 27  
 Hong, 477  
 Horn & Johnson, 201, 489–491  
 horn, flared, 81  
 Householder, 521  
 Hu, 319, 320  
 Huang  
     Feng, 301  
     Hong-Xuan, 360  
 hull, 50  
     affine, 33, 50–53, 65, 118, 634  
     cone, 51, 111, 125, 138, 149–151, 153, 573  
     cone, EDM, 376, 441  
     cone, positive semidefinite, 51  
     correlation matrices rank-one, 51  
     empty set, 51  
     point, 34, 51  
     unique, 50  
     conic, 52, 58, 88, 118, 145  
     empty set, 58  
     convex, 50–52, 88, 118, 342  
         cardinality-one vectors, 175, 281  
         cone, 88, 101, 130  
         dyads, 55, 56  
         empty set, 53  
         extreme directions, 87  
         orthogonal matrices, 58  
         orthonormal matrices, 58  
         outer product, 53–57, 248, 533  
         permutation matrices, 57  
         positive semidefinite cone, 101  
         projection matrices, 53, 248, 533  
         rank-one matrices, 55, 56  
         rank-one symmetric matrices, 53, 101  
         unique, 51  
 Huo, 273  
 hyperboloid, 514  
 hyperbox, 67, 603  
 hypercube, 44, 67, 175, 603  
     nonnegative, 237  
     slice, nonnegative, 273, 307  
 hyperdimensional, 410, 547, 549  
 hyperdisc, 552  
 hyperparallelepiped, 118  
 hyperplane, 58, 60–63, 194  
     angle, 69  
     hypersphere radius, 62  
     independence, 70  
     intersection, 70, 138  
         cone, 116  
         convex set, 65, 435  
         polyhedra, 116  
     movement, 62, 72, 73, 177  
     normal, 61  
     nullspace, 60, 61  
     parallel, 33  
     separating, 68, 130, 132, 133  
     strictly, 68  
     supporting, 65–67, 159, 197–199  
         cone, 122, 123, 138, 145  
         cone, EDM, 441  
         cone, PSD, 97, 98  
         exposed face, 75  
         polarity, 67  
         strictly, 67, 82, 197, 198  
         trivially, 65, 140  
         unique, 197, 198  
     vertex-description of, *see* description  
 hypersphere, 53, 260, 350, 389  
     circum-, 321, 381  
     Eternity II, 321  
     intersection, 289  
     packing, 354  
     radius, 62  
 hypograph, 186, 191, 218, 565  
 Hz, 630  
     – *I* –, 632  
 idempotent, *see* matrix  
 Identity, *see* matrix  
 iff, 633  
 image, 40, 41  
     affine, 40  
     inverse, 40, 41, 143  
     affine, 40  
         cone, 93, 130, 144, 187, 204, 209, 228  
         magnetic resonance, 300  
 in, 630  
 inactive, *see* active set  
 independence  
     affine, 59, 63, 64, 112, 115  
     preservation, 64  
     subset, 64, 70, 71  
     conic, 22, 59, 111–115, 238, 320  
         *versus* dimension, 111  
         dual cone, 153  
         extreme direction, 114  
         preservation, 112  
         rows, 153, 157, 160, 398, 399  
         unique, 111, 115, 150, 151  
     hyperplane, 70  
     linear, 32, 59, 111, 112, 114, 238, 282, 518  
     matrix, 494  
     of subspace, 32, 624  
     preservation, 32  
 inequality  
     active, 140, 244, 631  
     angle, 344, 409  
         matrix, 369  
     Bessel, 576

- calculus, 7  
 Cauchy-Schwarz, 71, 610  
 constraint, 161, 244  
 determinant, 496  
 diagonal, 496  
 eigenvalue, 493, 537  
 generalized, 84, 122, 131  
     dual, 131  
     dual PSD, 141  
     partial order, 84, 91  
 Hardy-Littlewood-Pólya, 463  
 Identity, 492, 532  
 inverse, 181, 496  
 Jensen, 197  
 linear, 20, 31, 60, 132  
     matrix, 141, 143, 144, 228, 231  
 log, 191, 565  
 Löwner, 91, 492, 495, 496  
     inverse, 496  
 norm, 235, 635  
     triangle, 173  
 rank, 496  
 semidefinite, *see* Löwner  
 singular value, 539  
 spectral, 396  
 sum, 147  
 trace, 495, 496  
 triangle, 342, 343, 368, 383–387, 392, 407–411, 418, 419  
     norm, 173  
     strict, 385  
     unique, 408  
     vector, 173  
 variation, 135  
 volume, 411  
 inertia, 395, 499  
     complementary, 397, 499, 500  
     preservation, 395, 494  
     Sylvester’s law, 494  
 infimum, 529, 633, *see* minimum  
     of concave functions, 219  
     of quasiconcave functions, 219  
 inflection, 200  
 injection, 631, *see* injective  
 injective, 44, 45, 85, 166, 373, 631  
     Gram-form, 371, 375, 376  
     inner-product form, 377, 378  
     invertible, 44, 567, 568  
     non, 46  
     nullspace, 44  
 innovation, rate, 299  
 interior, 34, 35, 39, 76, 633  
     0, 82, 85  
     algebra, 34, 39, 143  
     empty, 34  
     of affine, 34  
     of cone, 143–145  
         EDM, 416, 422  
         membership relation, 131, 152  
         polyhedral, 145  
         PSD, 90, 230  
     of halfspace, 39  
     of point, 34, 35  
     of ray, 77  
     -point, 34  
 antisymmetry, 116  
 barrier, 8, 223, 354  
 complementarity, 223  
 dimension, 309, 356, 476  
 intensity, 223, 455  
 method, 221, 236, 266  
 rank, 449  
 relative, 34, 39, 76  
 transformation, 143  
 intersection, 39  
     affine, 579  
     cone, 81  
         dual of, 130  
         of dual, 130  
         pointed, 82  
     epigraph, 185, 191, 219  
     face, 138  
     halfspace, 61, 129, 141  
         cone, 116  
         dual cone, 122  
         polyhedra, 116  
     hyperplane, 70, 138  
         cone, 116  
         convex set, 65, 435  
         polyhedra, 116  
     line with boundary, 36  
     nullspace, 73  
     planes, 71  
     positive semidefinite cone  
         affine, 110, 231  
         geometric center subspace, 513  
         line, 359  
         subspace, 73  
         tangential, 37  
 invariance  
     closure, 34, 116, 143  
     cone, 81, 112, 130, 416  
         pointed, 82, 85, 119  
     convexity, 39, 40  
     dimension, 45  
     EDM, 370  
         Gram-form, 371  
         inner-product form, 373  
         reflection, 371  
         rotation, 371  
         scaling, 347  
         translation, 370, 371, 375, 390  
     ellipsoid, 42  
     function, convex, 170, 185, 218  
     isometric, 21, 360, 370, 373, 400, 401  
     optimization, 173, 184, 248, 529, 530, 631  
     orthogonal, 45, 406, 458, 525, 527  
     rotation, 102  
     scaling, 218, 529  
     set, 394  
     translation  
         function, 170, 218, 529  
         subspace, 439, 593, 628  
 inverse  
     image, *see* image  
     matrix, *see* matrix  
     minimization, *see* function, fractional  
     nullspace, 41, 44, 568  
     pseudo, *see* matrix, pseudoinverse  
 inversion

- conic coordinates, 166
- Gram-form, 376
- invertible
  - injective, 44, 567, 568
  - operator, 42, 44, 45, 143
- is, 624
- isedm(), 344, 349, 460
- Ising, 324, 325
- isometry, 45, 406, 459
- isomorphic, 43, 46, 48, 623
  - face, 97, 351, 424
  - isometrically, 36, 44, 48, 94, 141
  - non, 513
- isomorphism, 43, 143, 144, 376, 378, 406, 416
  - isometric, 44, 45, 47, 49, 406
    - symmetric hollow subspace, 49
    - symmetric matrix subspace, 47
    - projection, 44
  - isotonic, 404
- iterate, 606, 614, 617
- iteration
  - alternating projection, *see* projection
  - convex, *see* convex
- $\mathcal{J}$  –, 301
- Jacobian, 193
- Jensen, 197
- Johnson, 201, 489–491
- Jordan, 297
- $\mathcal{K}$  –, 77
- Kaczmarz, 579
- Karhunen-Loéve transform, 400
- Karmarkar, 223
- kernel, *see* nullspace
- Kimmel, 24
- Kirschhoffer, 307
- kissing, 129, 174, 175, 280–282, 570
  - number, 354
  - problem, 354
- KKT conditions, 161, 162, 461, 541
- Klanfer, 350
- Klee, 87
- Korkine, 356
- Krein-Rutman, 130
- Kreyszig, 342, 383
- Kronecker product, *see* product
- Kuhn, 352
- $\mathcal{L}$  –, 184
- Lagrange, 125, 235
  - multiplier, 135, 540
  - sign, 167
- Lagrangian, 167, 351, 541, 622
  - MAX CUT, 296
  - origins, 172
- Lanckriet, 297
- Lanczos, 531
- Laplace
  - transfer function, 209–211, 267, 268
  - zero implementation, 210
- transform, 209–211, 267, 268
- Lasserre, 230
- Lasso, 299
- lattice, 255–258, 261–263
  - regular, 27, 255
- Laurent, 239, 240, 392, 430, 434, 615
  - vertex, 388
- Lauterbur, 300
- Law, 300, 548
- least
  - 1-norm, 174, 175, 234, 280
  - nonnegative, 280, 281
  - energy, 350, 473
  - norm, 175, 235, 570
  - squares, 253, 254, 570
- Legendre-Fenchel transform, 472
- Lemaréchal, 19, 67
- Lewis, 74, 84, 116, 595
- Li, 32, 64, 111, 114, 503, 518, 623
- Liang, 360
- limit, 633
- line, 32, 33, 51, 194, 195, 204
  - cone, 81
  - fit, 299, 338, 339, 585, 586
  - tangential, 38
- linear
  - algebra, 69, 485
  - bijection, *see* bijective
  - complementarity, 161
  - function, *see* operator
  - independence, *see* independence
  - inequality, 31, 60, 132
    - matrix, 141, 143, 144, 228, 231
  - injective, *see* injective
  - map, *see* transformation
  - operator, 40, 48, 169, 183, 348, 376, 377, 406, 426, 485, 594
    - projector, 46, 574, 577, 579, 596
  - program, *see* program
  - regression, 299, 585, 586
  - surjective, *see* surjective
  - transformation, *see* transformation
- list, 19, 26, 50, 625, 627, 631
  - generating, 118, *see* generators
  - of points in  $X$ , 50
  - reconstruction, *see* reconstruction
- Littlewood, 463
- Liu, 415, 420, 426, 476
- localization, 21, 359
  - sensor network, 252–266
  - standardized test, 255
  - unique, 22, 253, 254, 358, 360
  - wireless, 252, 259, 364
- log, 564, 565
  - barrier, *see* interior-point
  - constraint, 191
  - convex, 218
  - det, 206, 363, 473, 556, 563
- Löwner inequality, 91, 492, 495, 496
  - inverse, 496
- LP, 630, *see* program, linear
- Luenberger, 596
- Lustig, 301
- Lyapunov, 487

- $\mathcal{M}$  –, 155
- machine
  - control, 183, 483, 616
  - learning, 21, 24, 483, 585
- Magaia, 509
- majorization, 488
  - cone, 488
  - symmetric hollow, 488
- manifold, 21, 25, 57, 58, 96, 436, 437, 458, 482, 525, 535
- map, *see* transformation
  - isotonic, 402, 404
  - USA, 26, 402–407
- Mardia, 460
- Markov process, 441
- Marsden, 342
- Marziliano, 299
- Mason flowgraph, 211
- mater, 544
- Mathar, 442, 524
- Mathematica, 211, 338
- Matlab, 28, 526
  - backslash, 234
  - cartography, 403, 405
  - conic independence, 111
  - convex iteration
    - accelerant, 339
    - stall, 287
  - cvx, 291
  - EDM, 460
  - Eternity II, 321
  - FAST MAX CUT, 297
  - Hadamard product, 548
  - Kronecker product, 547
  - MRI, 299, 302, 305–307
  - noise, 264
  - notation, 276
  - polynomial feasibility, 292
  - quantum, 327
  - signal dropout, 280
  - SVD, 337, 507
- matrix, 544
  - 0**, 316
  - adjacency, 320
  - angle
    - interpoint, 348, 350
    - relative, 369
  - antisymmetric, 47, 488, 525, 628
    - antihollow, 48, 49, 456, 628
    - subspace, 47
  - arrow, 528
  - auxiliary, 350, 522, 524
    - Householder, 522
    - orthonormal, 524
    - projector, 522
    - Schoenberg, 347, 523, 584, 626
    - table, 524
  - binary, *see* matrix, Boolean
  - Boolean, 51, 236, 238, 388, 389, 627
    - orthogonal, 57
  - bordered, 386, 395
    - arrow, 528
  - calculus, 543
  - circulant, 206, 505, 522
  - permutation, 302, 311, 315–317
  - symmetric, 302
- commutative, *see* product
- completion, *see* problem
- correlation, 388, 400
  - rank-one, 51, 389
- decomposition, *see* decomposition
- determinant, *see* determinant
- diagonal, *see* diagonal
- diagonalizable, *see* diagonalizable
- direction, *see* direction
- distance
  - 1-norm, 202, 394
  - absolute, 341, 394, 418, 419, 459, 468–470, 622
  - binary, 202
  - Euclidean, *see* EDM
  - doublet, 371, 519, 594
    - nullspace, 519, 520
    - range, 519, 520
  - dyad, 514, 515
    - decomposition, 514
    - hull, 53
  - independence, 73, 148, 401, 422, 517, 518
  - negative, 515
  - nullspace, 515, 516
  - projection on, 587, 589
  - projector, 516, 517, 572, 587
  - pseudoinverse, 516
    - range, 515, 516
    - sum, 503, 506, 518, 519
    - symmetric, 101, 497, 498, 517
  - EDM, *see* EDM
  - elementary, 394, 520, 522
    - nullspace, 520, 521
      - range, 520, 521
    - entry, *see* entry
    - Euclidean distance, *see* EDM
    - exponential, 206, 565
    - factorization, *see* decomposition
    - fat, *see* wide
    - Fourier, 45, 301
    - fractional, *see* function
    - full-rank, 70
    - Gale, 360
    - geometric centering, 350, 371, 416, 439, 522
    - Gram, 259, 348, 357
      - unique, 371, 375
    - Hermitian, 488
    - hollow, 48, 49, 622, 628
    - Householder, 521, 526
      - auxiliary, 522
    - idempotent, 574
      - nonsymmetric, 571
      - range, 571
      - symmetric, 575, 577
      - transpose, 571
    - Identity, 632
      - direction vector, 250, 330
      - eigen, 492
      - Fantope, 53, 54
      - inequality, 492, 532
      - orthogonal, 57, 289, 521, 525
      - permutation, 57, 289, 521
      - positive definite, 57, 101, 102, 498, 525

- projection, 498
- indefinite, 329, 395
- indices, 342
- inverse, 205, 504, 556, 621
  - injective, 568
  - minimization, *see* function, fractional
  - positive definite, 496
  - product, 487
  - symmetric, 505
  - transpose, 487, 621
  - update, 517
- Jordan form, 492, 504
- measurement, 455
- nonexpansive, 525
- nonnegative, 395, 427, 446, 455, 602
  - definite, *see* matrix, positive semidefinite
  - factorization, 251
- nonsingular, 504, 505
- normal, 44, 492, 505, 506
- normalization, 237
- nullspace, *see* nullspace
- orthogonal, 45, 372, 505, 507, 525–527
  - manifold, 525
  - permutation, 57, 288, 289, 525
  - product, *see* product
  - symmetric, 526
  - transpose, 525
- orthonormal, 45, 53, 58, 248, 287, 371, 506, 507, 524, 533, 576, 577
  - manifold, 58
  - nonexpansive, 525
  - pseudoinverse, 568
  - square, 525
- partitioned, 499–502, *see* Schur
- permutation, 57, 288, 311, 314, 382, 525
  - $\delta$ , 487, 505
  - circulant, 302, 311, 315–317
  - extreme point, 57, 316
  - orthogonal, 57, 288, 289, 525
  - positive semidefinite, 57
  - product, *see* product
  - symmetric, 302, 521, 537
- positive definite, 91
  - eigenvalues, 492, 496, 509
  - inverse, 496
  - positive real numbers, 201
  - singular values, 509
- positive factorization, 251
- positive semidefinite, 90–110, 141, 395, 488–502
  - completely, 81, 251
  - difference, 110
  - eigenvalues, 492, 496, 509
  - extreme direction decomposition, 102
  - nonnegative *versus*, 395
  - nonsymmetric, 491
  - permutation, 57
  - projection, 498, 576
  - pseudoinverse, 569
  - rank, *see* rank
  - singular values, 509
  - square root, 496, 505, 622
  - sum, eigenvalues, 497
  - sum, nullspace, 73
  - sum, rank, *see* rank
- symmetry *versus*, 489
- test-domain, 488, 591
- zero entry, 510
- product, *see* product
- projection, 522, 571, 576
  - diagonal, 94
  - eigenvalue, 572, 574, 576
  - nonorthogonal, 571
  - orthogonal, 575
  - positive semidefinite, 498, 576
- product, *see* product
- rank, *see* rank
- transpose, 571
- pseudoinverse, 41, 94, 149, 194, 235, 499, 567–569
  - 0, 569
  - by SVD, 510
  - inverse, 568, 569
  - nullspace, 41, 567, 568
  - of dyad, 516
  - of projection, 568, 576
  - of vector, 568
  - orthonormal, 568
  - positive semidefinite, 569
- product, *see* product
- range, 41, 567, 568
  - symmetric, 510, 569
  - transpose, 147
  - unique, 194, 567
- quotient, *see* Hadamard
- range, *see* range
- rank, *see* rank
- reflection, 371, 372, 526, 537
  - of range, 527
  - rotation, 371–373, 525–527, 536
    - of range, 526, 527
- Schur-form, *see* Schur
- similarity, 496, 590
  - sign, 498
  - simple, 515
  - skinny, *see* thin
  - sort index, 403
  - square, 205
  - square root, 496, 622, *see* matrix, positive semidefinite
- Stiefel, *see* matrix, orthonormal
- stochastic, 57, 288
- sum
  - eigenvalues, 497
  - nullspace, 73
  - rank, *see* rank
- symmetric, 42, 46–49, 488, 505, 628
  - antihollow, 48, 628
  - eigenvalues, 505, 509
  - Hermitian, 488
  - inverse, 505
  - pseudoinverse, 510, 569
  - real numbers, 201
  - singular values, 509
  - subspace, 46
  - symmetrized, 489
- thin, 45, 149, 622
- trace, *see* trace
- transpose, 621
  - conjugate, 488, 525, 621

- idempotent, 571
- inverse, 487, 621
- unitary, 525
  - symmetric, 301, 303
  - wide, 46, 70, 329, 622
  - zero definite, 513, 514
- max, 633
- MAX CUT problem, 294
- maximal, 633
- maximization, *see also* supremum
  - Eternity II, 321
  - of distance, 439, 597
  - of norm, 322, 530
  - of Procrustes, 537
  - of quadratic, 296
    - nonconcave, 322, 532, 541
  - of Rayleigh quotient, 532
  - of trace, 439, 530
- maximize, 631, 633
- maximum, 633
  - of totally ordered set, 633
  - variance unfolding, 436, 441
- McCormick, 223
- membership, 51, 84, 91, 634
  - boundary, 138, 231
  - relation, 131, 152
    - affine, 172
    - boundary, 132, 138, 231
    - discretized, 136, 137, 141, 450
    - EDM, 450
    - interior, 131, 152, 229
    - normal cone, 158
    - orthant, 51, 137, 138
    - subspace, 152
- Menger, 382, 396, 410, 412, 463, 480
- metric, 45, 202, 341
  - audio, 215
  - postulate, 342
  - property, 342, 383
    - fifth, 343–345, 407–409, 413
  - space, 342
- min, 633
- minimal, 633
  - cardinality, *see* minimization
  - element, 83, 84, 171, 172
  - minimum, 171, 172
  - rank, *see* minimization
  - set, *see* set
- minimization, 134, 158, 254, 529, 552
  - affine function, 184
  - fractional, *see* function, fractional
  - norm
    - 1-, 174, *see* problem
    - 2-, 173, 352, 586
    - $\infty$ -, 174
  - Frobenius, 173, 189, 190, 193, 194, 459, 467, 592
  - nuclear, 330, 472, 531
- of cardinality, 174, 175, 271–275, 285, 286, 304
  - Boolean, 234
  - by perturbation, 244
  - nonnegative, 176–179, 271–275, 280–284
  - rank connection, 249
  - rank, full, 238
  - rank-one, 297
- reduction, 244–246
- of distance, 575, 592
- of quadratic
  - binary, 326, 327
  - convex, 162, 190, 406, 540
  - nonconvex, 532, 540
- of rank, 247–251, 287, 332, 333, 474, 480
  - by perturbation, 239
  - cardinality connection, 249
  - cardinality constrained, 297
  - indefinite, 329
  - reduction, 224, 238–242
  - wide, 329
- of Rayleigh quotient, 532
- of trace, 249, 250, 330, 357, 471–473, 500, 531
  - on hypercube, 67
- minimize, 631, 633
- minimizer, 171
  - unique, 170, 171
- minimum, 633
  - element, 83, 84, 171, 172
    - unique, 83, 84, 171, 172
  - global, 169, 171, 188, 206, 552, 593
  - local, 171
    - minimal, 171, 172
    - norm, *see* least
    - of convex function, 171
    - of totally ordered set, 633
    - phase, 209, 267
      - unique, 169, 552
  - Minkowski, 40, 117
  - MIT, 331
  - Mizukoshi, 50
  - modulus, 633
  - molecular conformation, 21, 23, 27, 366
  - Monckton, 307, 309, 310, 316, 321
  - monotonic, 189, 195, 251, 273, 611, 613
    - Fejér, 610
    - function, *see* function
    - gradient, 196
    - noisily, 251, 273
  - Moore-Penrose
    - conditions, 567, 569
    - inverse, *see* matrix, pseudoinverse
  - Moreau, 161, 600
  - Morrison, 517
  - Motzkin, 109, 134, 292, 595
  - Mount, 215
  - MRI, 299–303, 627
  - Muller, 509
  - multidimensional
    - function, *see* function
    - objective, 83, 84, 171, 172, 270, 304
      - scaling, 21, 399, 419, 460, 484
        - ordinal, 404
    - multilateration, 253, 364
  - multiobjective optimization, 83, 84, 171, 172, 270, 304
  - multipath, 364
  - multiplicity, 493, 504
  - Musin, 354, 356
  - $-\mathcal{N}-$ , 69
  - necessary, 624

- neighbor  
 base station, 365  
 nearest, 25, 436–440  
 pixel, 301, 304, 305  
 qubit, 323, 324  
 neighborhood graph, 436, 437  
 Németh, 135, 161, 604  
 Nemirovski, 455  
 nested  
   sequence, 385  
   sublevel set, 199  
 Newton, 354, 543  
   direction, 193  
 Nigam, 366  
 Nirenberg, 596  
 node, 255–258, 261–265, 294, 295  
 nondegeneracy, 342  
 nonexpansive, 161, 525, 571, 576, 602  
 nonnegative, 382, 635  
   constraint, 57, 71, 72, 161, 173–177, 222, 286  
   factorization, 251  
   part, 180, 456, 603, 604, 625  
   polynomial, 85, 500  
 nonnegativity, 173, 342  
 nonorthogonal  
   basis, 149, 580  
   projection on, 580  
   projection, *see* projection  
 nonvertical, 194  
 norm, 171, 173–176, 271, 635  
   0, 510  
   0-, *see* 0-norm  
   1-, *see* 1-norm  
   2-, *see* 2-norm  
    $\infty$ -, *see*  $\infty$ -norm  
   k-largest, 178, 179, 274, 321, 635  
     gradient, 180  
     monotonicity, 179  
    $\ell$ -, 635  
   ball, *see* ball  
   constraint, 286, 289, 291, 293, 510, 585  
     Frobenius, 586  
     spectral, 58  
   convexity, 58, 170–175, 191, 196, 346, 530, 635  
   dual, 125, 530, 635  
   equivalence, 173  
   Euclidean, *see* 2-norm  
   Frobenius, 44, 171, 249, 458, 635  
     constraint, 586  
     maximal, 322  
     minimal, 173, 189, 190, 193, 194, 467, 592  
     Schur-form, 190  
     gradient, 179, 180, 192, 559, 561, 575, 583,  
       584, 592, 595  
   homogeneity, 173  
   inequality, 235, 635  
     triangle, 173  
   Ky Fan, 530  
   least, 175, 235, 570  
   nuclear, *see* nuclear  
   of difference, 346  
   of dyad, 516  
   of outer product, 516  
   orthogonally invariant, 45, 406, 458, 525, 527  
   property, 173  
 regularization, 294, 570  
 residual, 298  
 spectral, 58, 191, 249, 458, 468, 509, 516, 521,  
   530, 604, 635  
 ball, 58, 604  
 constraint, 58  
 dual, 530, 635  
 inequality, 635  
 Schur-form, 191, 530  
 square, 171, 196, 561, 583, 584, 592, 635  
   vs. square root, 173, 190, 196, 458, 575, 586  
 normal, 61, 595  
   cone, *see* cone  
   equation, 569  
   facet, 145  
   gradient, 159, 193  
   inward, 60, 441  
   outward, 60  
   vector, 61, 595, 619  
 not, 624  
   -necessarily, 633  
 Notation, 621  
 NP-hard, 297, 307, 475  
 nuclear  
   magnetic resonance, 23, 366  
   norm, 249, 330, 472, 530, 531, 635  
     ball, 55–57, 531  
 nullspace, 41, 61, 69–74, 522, 568, 594, 627  
   algebra, 515  
   basis, 61, 70, 73, 74, 100, 135, 151, 160, 507,  
     575, 583, 584, 628  
     Schur-form, 500  
     vectorized, 74, 100  
   cone, pointed, 119  
   decomposition, 61, 570  
   dimension, 69  
   doublet, 519, 520  
   dyad, 515, 516  
   elementary matrix, 520, 521  
   -form, 69  
   hyperplane, 60, 61  
   injective, 44  
   intersection, 73  
   inverse, 41, 44, 568  
   of product, 380, 494, 515, 519, 576  
   operator, *see* operator  
   orthogonal complement, 69  
   projector, 575  
   pseudoinverse, 41, 567, 568  
   set, 138  
   sum, 73  
 numerical precision, *see* precision  
  
   –  $O$  –, 632  
 objective, 67, 68, 127, 134, 158, 160, 631  
 convex, 221  
   strictly, 469, 473  
   linear, 183, 261, 357  
   multidimensional, 83, 84, 171, 172, 270, 304  
   nonlinear, 184  
   polynomial, 291  
   quadratic  
     convex, 162, 190, 406, 477, 540  
     nonconvex, 322, 540

- real, 171
- value, 231, 241
- offset, *see* invariance, translation on, 630
- one-to-one, 631, *see* injective
- only if, 624
- onto, 631, *see* surjective
- op amp, 210, 211
- open, *see* set
- operator, 624, 632
  - adjoint, 42, 43, 130, 447, 621, 626
  - self, 42, 351, 446, 485
- affine, *see* function
- injective, *see* injective
- invertible, 42, 44, 45, 143
- linear, 40, 48, 169, 183, 348, 376, 377, 406, 426, 485, 594
  - projector, 46, 574, 577, 579, 596
- nonlinear, 170, 184
- nullspace, 371, 375–377, 594
- permutation, 178, 272, 462–465, 481, 486, 626
- projection, *see* projector
- quadratic, convex, 347, 370
- surjective, *see* surjective
- unitary, 45, 525
- optimal, 7, 19, 631
  - optimum, 171, 172, 631
- optimality
  - condition, 232
    - conic problem, 127, 160, 162
    - direction vector, 248, 272
    - directional derivative, 552
    - dual, 125
    - first-order, 134, 135, 158–162
    - KKT, 161, 162, 461, 541
    - linear program, 127
    - semidefinite program, 233
    - Slater, 127, 167, 229, 232, 233
    - unconstrained, 135, 193, 194, 592
  - constraint
    - conic, 160, 162
    - equality, 135, 162
    - inequality, 161
    - nonnegative, 161
  - global, 7, 19, 171, 188, 221
    - convex iteration, 248, 250, 251, 272, 336
  - local, 7, 19, 171, 221
    - convex iteration, 251, 273, 276, 285, 336
- optimization, 19, 221
  - algebra, 173, 529, 530
  - combinatorial, 57, 68, 177, 234, 292, 294, 311, 315, 326, 389, 483
  - Procrustes, 288
  - conic, *see* problem
  - convex, 7, 19, 171, 221, 357
    - art, 8, 238, 254, 484
    - fundamental, 221
    - solution set, 171, 333
  - invariance, 173, 184, 248, 529, 530, 631
  - multiobjective, 83, 84, 171, 172, 270, 304
  - programming, 221
  - quantum, 324–326
  - tractable, 221
  - vector, 83, 84, 171, 172, 270, 304
- optimum, 631
- optimal, 171, 172, 631
- order
  - natural, 42, 485, 632
  - nonincreasing, 464, 496, 505, 506, 626
  - O, 632
  - of projection, 456–458, 605, 617
  - partial, 51, 81, 83, 84, 91, 137, 138, 145, 146, 495, 632, 634
  - generalized inequality, 84, 91
  - orthant, 138
  - property, 84
  - total, 83, 84, 122, 633
  - transitivity, 84, 91, 495, 632
- origin, 31, 32, *see* zero
- cone, *see* cone
- projection of, 63, 175, 570, 583
  - 1-norm, 174
  - projection on, 599, 604
  - subspace, 32, 47, 627, 628, 632
  - translation to, 379
- Orion nebula, 21
- orthant, 33, 120, 122, 131, 138, 149, 158, 631
  - dual, 141
  - extreme direction, 72
  - nonnegative, 132, 446
- orthogonal, 42, 575, 592, 628
  - basis, 48, 49, 73, 577
  - projection on, 580
  - complement, *see* complement
  - condition, 575
  - constraint, 335, 367, *see* Procrustes equivalence, 527
  - expansion, 48, 122, 148, 576, 577
  - invariance, 45, 406, 458, 525, 527
  - matrix, *see* matrix
  - projection, *see* projection
  - set, 138, 628
  - sum, 624
  - vector, 42, 628
- orthogonality, 575, 631
  - matrix, 628
  - set, 138, 628
  - vector, 42, 628
- orthonormal, 42
  - condition, 148, 149, 576
  - constraint, 287
  - decomposition, 576, 577
  - matrix, *see* matrix
- over, 630
- overdetermined, 622
- $-P-$ , 571
- PageRank, 441
- parallelepiped, 118
- parallelogram, 118
- Pardalos, 360
- Parhizkar, 20
- pattern recognition, 21
- Penrose conditions, 567, 569
- pentahedron, 411
- pentatope, 120, 411
- permutation, *see* matrix
  - constraint, *see* polyhedron, permutation
- perpendicular, *see* orthogonal

- Perron, 395  
 perturbation  
     cardinality, 244, 245  
     rank, 239–241  
 Pfender & Ziegler, 354  
 phantom, 300  
 phase, 634  
     minimum, 209, 267  
     -transition, 176, 275  
 Photoshop, 330  
 plan, 67  
 plane, 58  
     segment, 86  
 point  
     boundary, 78  
         of hypercube slice, 307  
     exposed, *see* exposed  
     extreme, *see* extreme  
     feasible, *see* solution  
     fixed, 251, 305, 594, 608–611  
     inflection, 200  
     -list, *see* list  
     minimal, *see* minimal element  
     vector, 87, 631  
 Pólya, 463  
 Polyak, 529  
 polychoron, 39, 116, 120, 410, 411  
 polyhedron, 36, 50, 67, 116–120, 122, 387  
     bounded, 53, 57, 117, 177  
     face, 115, 116  
     halfspace, *see* halfspace  
         -description, 116  
     norm ball, 177  
     permutation, 57, 58, 288, 290, 316, 326  
         vertex, 321  
     range form, 118  
     stochastic, *see* polyhedron, permutation  
     transformation  
         linear, 112, 116, 132  
         unbounded, 33, 36, 117  
         vertex, 57, 117, 177  
         -description, 53, 118  
 polynomial  
     constraint, 180, 181, 291  
     convex, 200  
     Motzkin, 292  
     nonconvex, 291  
     nonnegative, 85, 500  
     objective, 291  
     quadratic, *see* function  
 polytope, 116  
 positive, 635  
     completely, 81, 251  
     factorization, 251  
     semidefinite, *see* matrix  
     strictly, 383  
 power, *see* function, fractional  
 Précis, 381  
 precision  
     antisymmetry, 116  
     eigenvalue, 503  
     inequalities, complementary, 116  
     numerical, 216, 336, 337  
     perturbation, 240  
     quadruple, 216  
 solver, 304  
     barrier, 8, 223, 354  
     crippling, 223  
     dimension, 223, 455  
     hollowness, 483  
     interior-point, 236, 266, 309, 356, 476  
     variables, 223, 309  
 presolver, 139, 282–285, 309, 318–320  
     *en masse*, 139  
     aggregation, 284  
 primal  
     feasible set, 224, 228  
     problem, 127, 222  
         via dual, 125, 128, 129, 234  
 principal  
     component analysis, 297, 400, 484, 585, 586  
     eigenvalue, 297, 532  
     eigenvector, 297, 298, 532, 586  
     submatrix, 357, 408, 411, 425  
         face, 98  
         leading, 99, 383, 385  
         positive semidefinite, 497  
         rank, 98, 497  
 principle  
     halfspace-description, 60  
     separating hyperplane, 68  
     supporting hyperplane, 65  
 problem  
     1-norm, 174–177, 234  
         nonnegative, 176, 177, 272, 273, 275,  
             280–284, 472  
         signed, 285, 286  
     ball packing, 354  
     Boolean, 57, 68, 234, 292, 315, 326, 389  
     combinatorial, 57, 68, 177, 234, 292, 294, 311,  
         315, 326, 389, 483  
     Procrustes, 288  
     complementarity, 161  
         linear, 161  
     completion, 255–258, 264, 302, 309, 343, 344,  
         358, 360, 361, 386, 392, 413, 435–437,  
         484  
     geometry, 436  
     semidefinite, 614, 615  
     compressed sensing, *see* compressed  
     concave, 221  
     conic, 127, 160, 162, 172, 222  
         dual, 127, 222  
     convex, 127, 135, 160, 162, 174, 183, 188, 190,  
         223, 228, 466, 595, 631  
         definition, 221  
         geometry, 19, 23  
         nonlinear, 221  
         statement as solution, 8, 221  
         tractable, 221  
     dropout, signal, 276–280  
     dual, *see* dual  
     epigraph form, 189, 190, 478, 586  
     equivalent, 173, 184, 248, 260, 529, 530, 631  
     Eternity II, *see* Eternity  
     feasibility, 20, 111, 119, 138, 139, 632, 633  
         semidefinite, 231, 247, 248, 334  
     Gram-form, *see* problem, proximity  
     kissing, 354  
     MAX CUT, 294

- minimax, 51, 125, 128
- nonconvex, 8
  - EDM, PSD, 458, 464–466
  - LP, 71
  - map, 26, 402
  - polynomial, 291
  - projection, 466
  - stress, 475
  - nonlinear, 357
    - convex, 221
  - open, 26, 107, 251, 399, 426, 483
  - permutation, 57, 288, 290, 311, 314–316, 326
  - primal, *see* primal
  - Procrustes, *see* Procrustes
  - proximity, 455, 456, 458, 460, 468, 476
    - EDM, nonconvex, 464
    - Gram-form, 467, 470, 478
    - in spectral norm, 468
    - rank heuristic, 472, 474
    - semidefinite, 462, 469, 470, 477, 478
  - quadratic
    - binary, 326, 327
    - convex, 162, 406, 540
    - nonconvex, 322, 540
  - quasiconvex, 206, 214, 216
  - same, 631
  - sphere packing, 354
  - sstress, 458, 459
  - strain, 459
  - stress, 201, 458, 459, 475
  - tractable, 221
  - tug of war, 352
- procedure
  - alternating projection, 617
  - bisection, 212
  - cardinality reduction, 246
  - rank reduction, 239
- Procrustes
  - combinatorial, 288
  - compressed sensing, 290
  - diagonal, 540
  - linear program, 538
  - maximization, 537
  - orthogonal, 535
    - two sided, 537, 538
  - orthonormal, 287
  - permutation, 288, 535
  - symmetric, 539
  - translation, 536
  - unique solution, 535
  - vector, 288, 535
- product, 208, 547, 626
  - Cartesian, 40, 228, 612, 627
    - cone, 81, 82, 128
    - function, 172
  - commutative, 206, 233, 493, 494, 512, 513, 606
  - determinant, 494
  - eigenvalue, 493–495, 535
    - of, 494, 495, 509
  - empty, 626
  - function, 187, 197, 208, 209, 216
  - gradient, 546
  - Hadamard, *see* Hadamard
    - inner
      - EDM definition, 367
  - matrix, 348
  - vector, 42, 60, 153, 209, 223, 299, 367, 368, 525, 587, 625
  - vector, positive semidefinite, 141
  - vectorized matrix, 42–44, 492
  - zero, *see* product, zero
- inverse, 487
  - transpose, 487
- Kronecker, 94, 485–487, 538, 547, 548
  - $\delta$ , 487
  - determinant, 487
  - diagonal, 487
  - dimension, 547
  - eigen, 487
  - gradient, 560
  - inverse, 302, 487
  - of vectors, 548
  - orthogonal, 527
  - permutation, 486, 525
  - positive semidefinite, 346, 495
  - projector, 592
  - pseudoinverse, 487, 568
  - range, 592
  - rank, 487
  - trace, 487
  - transpose, 487
    - vectorization, 487
  - nullspace, 380, 494, 515, 519, 576
  - orthogonal, 525
  - outer, 53–57, 248, 533
    - norm, 516
    - positive semidefinite, 96
    - vector, *see* matrix, dyad
  - permutation, 525
  - positive definite, 490
    - nonsymmetric, 490
  - positive semidefinite, 495, 498
    - matrix outer, 96
    - nonsymmetric, 494
    - vector inner, 141
    - zero, 512
  - projection, 443, 457, 458, 604–606, 610
  - projector, 443, 587, 592, 606
  - pseudofractional, 186
  - pseudoinverse, 568
  - quasiconcave, 209
  - range, 494, 515, 519, 575
  - rank, 380, 493, 494
  - singular value of, 509
  - symmetric, 494
  - tensor, 547
  - trace, *see* trace
  - transpose, 487
    - inverse, 487
  - zero, 42, 61, 63, 73, 132, 512–514, 628
  - PSD, 91, 97–101, 233, 512, 513
- program, 67, 221
  - binary, 315, 317, 318
  - quadratic, 326
  - class, 224
  - convex, *see* problem
  - geometric, 7, 224
  - integer, 57, 68, 389
  - linear, 67, 71–73, 128, 173, 183, 221–224, 354, 632

- dual, 128, 222
  - prototypical, 128, 222, 632
- nonlinear
  - convex, 221
- plan, 67
- quadratic, 162, 223, 224, 406
  - binary, 326
  - nonconvex, 322
- quadratically constrained, 190, 224, 291
- second-order cone, 223, 224
- semidefinite, 128, 173, 183, 221–224, 290, 632
  - dual, 128, 222
  - equality constraint, 91
  - intensity, 223
  - prototypical, 128, 222, 232, 244, 632
  - Schur-form, 181, 189–191, 406, 469, 478, 530, 586, 602
  - type, 224
- programming, optimization, 221
- projection, 567, 571, 576
  - 1-norm, 604
  - $I - P$ , 574, 600, 601
- algebra, 579
- alternating, 354, 606, 607
  - angle, 611
  - convergence, 610, 612, 613, 615
  - distance, 608, 609
  - Dykstra algorithm, 606, 616
  - feasibility, 608–610
  - iteration, 478, 606, 611, 617
  - on affine  $\cap$  PSD cone, 613
  - on EDM cone, 477
  - on halfspaces, 608
  - on orthant  $\cap$  hyperplane, 611
  - on subspace/affine, orthogonal, 606
  - optimization, 608, 609, 616
  - over/under, 612
- biorthogonal expansion, 577, 580–582
- coefficient, 46, 431, 433, 498, 575, 588–591
  - nonorthogonal, 147, 587, 588
- cyclic, 607
- direction, 575, 577, 587, 588
  - nonorthogonal, 572, 573, 582, 587
  - parallel, 578, 579, 581, 582
- dual, *see* dual
- easy, 603
- eigenvalue, 588–590
- Euclidean, 108, 456, 458, 460, 462, 468, 476, 595
- isomorphism, 44
- matrix, *see* matrix
- minimum-distance, 135, 573, 575, 592, 593, 595, 596
- nonorthogonal, 349, 459, 468, 571, 573, 581, 587
  - of origin, 174, 175
- oblique, *see* projection, nonorthogonal
- of convex set, 42
- of hypercube, 44
- of matrix, 592, 593
- of origin, 63, 570, 583
  - 1-norm, 174, 175
- of polyhedron, 116
- of PSD cone, 95
- on 1-norm ball, 604
- on affine, 42, 579, 584
- hyperplane, 63, 583
- of origin, 63, 174, 175, 570, 583
- orthogonal, 606
- vertex-description, 584
- on basis
  - nonorthogonal, 580
  - orthogonal, 580
- on box, 603
- on cardinality  $k$ , 603
- on complement
  - algebraic, 574, 600
  - orthogonal, 600
- on cone, 599, 600, 604
  - dual, *see* dual
  - Lorentz (second-order), 604
  - monotone nonnegative, 398, 406, 463, 464, 604
  - polar, *see* dual, projection, on cone
  - polyhedral, 604
  - simplicial, 604
  - truncated, 602
- on cone, EDM, 468, 469, 476–480
  - boundary, 480
  - dual, *see* dual
- on cone, PSD, 108, 460–462, 604
  - boundary, 27, 466
  - cardinality constrained, 297
  - geometric center subspace, 443
  - rank constrained, 193, 462, 466, 532
  - rank-one, 297
- on convex set, 135, 595, 596
  - boundary, 466, 598
  - dual, *see* dual
  - in affine subset, 606
    - in subspace, 605
  - on convex sets, 607
- on dual, *see* dual
- on dyad, 587, 589
- on ellipsoid boundary, 286, 287
- on Euclidean ball, 603
- on function domain, 186
- on geometric center subspace, 443, 593
- on halfspace, 583
- on hyperbox, 603
- on hypercube, 603
- on hyperplane, 583, *see* projection, on affine
- on hypersphere, 603
- on intersection, 606–617
- on line, 584
- on origin, 599, 604
- on orthant, 464, 602, 603
  - in subspace, 464
- on range, 569, 570, 575
- on rank constrained matrices, 604
- on rowspace, 46, 575
- on simplex, 604
- on slab, 583
- on spectral norm ball, 604
- on subspace, 42, 575
  - Cartesian, 603
  - elementary, 582
  - face, 98
  - hollow symmetric, 460, 603
- matrix, 592

- orthogonal, 606
- orthogonal complement, 77, 122, 599
- polyhedron, 116
- symmetric matrices, 603
- on vector
  - cardinality  $k$ , 603
  - nonorthogonal, 587
  - orthogonal, 588
  - on vectorized matrix, 587, 588
  - one-dimensional, 587–589
  - order of, 456–458, 605, 617
  - orthogonal, 569, 570, 575, 588
  - product, 443, 457, 458, 604–606, 610
  - range rowspace, 594
  - semidefiniteness test as, 591
  - spectral, 439, 463, 464
  - norm, 468, 604
    - unique, 464
  - successive, 607
  - two sided, 592–594
  - umbral, 42
  - unique, 573, 575, 592, 595, 596
  - vectorization, 430, 593
- projector, 42, 575, 596, 632
  - affine subset, 579
  - auxiliary matrix, 522
  - characteristic, 574, 577
    - nonorthogonal, 571
    - orthogonal, 576
  - commutative, 606
    - non, 607
  - complement, 574
  - convexity, 574, 577, 595, 596
  - direction, *see* projection
  - dyad, 516, 517, 572, 587
  - fractional function, 188
  - linear operator, 46, 574, 577, 579, 596
  - nonexpansive, 161, 571, 576, 602
  - nonorthogonal, 349, 573, 578, 587
    - affine subset, 571
  - nullspace, 575
  - orthogonal, 575–578
  - product, 443, 587, 592, 606
  - range, 575
  - rank trace, 574, 577
  - rowspace, 46, 575
  - semidefinite, 498, 576
  - subspace, fundamental, 575
  - unique, 575, 593
- proper
  - cone, *see* cone
  - subset, 623
  - subspace, 32, 48, 49, 125, 149, 153, 441, 600
- prototypical
  - complementarity problem, 161
  - compressed sensing, 280
  - coordinate system, 122
  - linear program, 128, 222, 632
  - SDP, 128, 222, 232, 244, 632
  - signal, 299
- PSD, *see* matrix, positive semidefinite
- pseudoinverse, *see* matrix
- puzzle, Eternity II, *see* Eternity
- pyramid, 412
- $-Q-$ , 525
- quadrant, 33, 631
- quadratic, *see* function
- quadrature, 310, 311, 372, 373
- quantum, 295, 324–326, 354
  - Eternity II, 326, 327
  - gap, 326–328
- quartix, 545, 631
- quasi-, *see* function
- qubit, 324
  - coupling, 324
  - neighboring, 323, 324
- QUBO, 326, 327, 630
- quotient
  - Hadamard, *see* Hadamard
  - Rayleigh, 591
    - optimization, 532
- $-\mathcal{R}-$ , 69
- range, 41, 52, 69, 568, 592, 627
  - basis, 87, 507, 575, 580, 627
    - complement, 151
    - vectorized, 74
  - dimension, 45
  - doublet, 519, 520
  - dyad, 515, 516
  - EDM, 422
  - elementary matrix, 520, 521
  - form, 69
    - polyhedron, 118
  - idempotent, 571
  - of product, 494, 515, 519, 575
    - Kronecker, 592
  - of vectorization, 587, 588, 590–592
  - orthogonal complement, 69
  - projector, 575
  - pseudoinverse, 41, 567, 568
  - reflection, 527
  - rotation, 526, 527
  - rowspace projection, 594
- rank, 493, 494, 506, 507, 634
  - $-\rho$  subset, 96, 109, 460–466, 480, 484
  - constraint, *see* constraint
- convex
  - envelope, *see* convex
  - iteration, *see* convex
  - subsets, 108
- diagonalizable, 493
- dimension, 111, 380, 381, 634
  - affine, 381, 474
  - EDM, 381, 422
  - full, 70, 238
  - heuristic, 472, 474
  - indefinite, 329
  - inequality, 496
  - Kronecker, *see* product
  - log det, 473, 474
  - minimization, *see* minimization
  - monotonicity, 196, 471
  - of product, 380, 493, 494
  - of symmetric, 98
  - of transpose, 380, 494
  - one, 498, 515

- Boolean, 51, 236, 238, 388, 389  
 convex iteration, 287, 333–336, 483  
 hull, *see* hull, convex  
 modification, 517, 520  
 PSD, 101, 497, 498, 517  
 subset, 193  
 symmetric, 101, 497, 498, 517  
 transformation, 334  
 update, *see* rank-one, modification  
 partitioned matrix, 501  
 positive semidefinite cone, 224, 496  
 convex subsets, 108  
 face, 97–100  
 Précis, 381  
 projection matrix, 574, 577  
 quasiconcavity, 107, 196, 471  
 reduction, *see* minimization  
 regularization, 270, 287, 294, 333, 474, 483  
 Schur-form, 501, 502  
 sum, 107  
 positive semidefinite, 107, 108, 493, 494  
 trace  
 gap, 471  
 heuristic, 249, 250, 330, 471–473, 531  
 zero, 97, 109, 225, 226, 250
- ray, 77  
 boundary of, 77  
 cone, 77, 81, 82, 112, 119  
 boundary, 85  
 EDM, 88  
 positive semidefinite, 87  
 extreme, 86, 428  
 interior, 77
- Rayleigh quotient, 591  
 optimization, 532
- realizable, 19, 353, 459, 460
- reconstruction, 399  
 isometric, 360, 403  
 isotonic, 402–407  
 list, 399, 403  
 unique, 344, 357–361, 370, 373, 436
- recursion, 305, 449, 485, 497
- reflection, 371, 372, 526, 537  
 invariance, 371  
 of matrix, 527  
 of range, 527  
 prevention, 372  
 signal, 20, 364  
 vector, 526
- reflexivity, 84, 495, 632
- regular  
 lattice, 27, 255  
 simplex, 120, 411  
 tetrahedron, 409, 419
- regularization  
 cardinality, 277, 298  
 norm, 294, 301, 304, 570  
 rank, 270, 287, 294, 333, 474, 483
- relative, 34, 633
- relaxation, 57, 68, 177, 235, 237, 359, 389, 405
- residual  
 amplitude, 176  
 norm, 298
- reweighting, 273, 473
- Riemann, 223
- Riordan, 309  
 Robberto, 21  
 robotics, 22, 28  
 Rockafellar, 67, 85, 86, 116, 125  
 room geometry, 20, 21  
 root, *see* function, fractional  
 rotation, 371, 525, 526, 536  
 invariance, 102, 371  
 of cone, *see* cone, rotated  
 of matrix, 527  
 of range, 526, 527  
 quadrature, 311, 372, 373  
 vector, 105, 372, 373, 525
- round, 633
- rowspace, 41, 69, 568, 627  
 projector, 46, 575  
 range projection, 594
- Rutman, 130
- $-\mathcal{S}-$ , 119
- saddle, 129  
 value, 127, 129, 223
- same, 631
- Saul, 25, 436, 438, 440
- Saunders, 216, 315, 317, 319, 517
- scaling, 399, *see also* homogeneity  
 frequency, 216  
 invariance, 218, 347, 529  
 multidimensional, 21, 399, 419, 460, 484  
 ordinal, 404  
 unidimensional, 474
- Schütte & van der Waerden, 354
- Schoenberg, 25, 341, 349, 387, 392, 394, 405, 459, 460  
 auxiliary matrix, 347, 523, 584, 626  
 criterion, 25, 26, 349, 350, 396, 416, 427, 430, 447, 450, 451, 453, 591
- Schur, 488  
 complement, 213, 260, 359, 499, 501  
 conditions, 499, 510, 528, 540  
 $\text{-form}$ , 499–502  
 anomaly, 190  
 constraint, 58  
 convex set, 93  
 norm, Frobenius, 190  
 norm, spectral, 191, 530  
 nullspace, 500  
 quadratic, 500  
 rank, 501, 502  
 semidefinite program, 181, 189–191, 406, 469, 478, 530, 586, 602  
 sparse, 501
- Schwarz, 71, 610
- SDP, *see* program, semidefinite
- semidefinite  
 positive, *see* matrix  
 program, *see* program
- sensor, 22, 252, 253, 255–258, 261–265  
 network localization, 252–266, 357, 363
- sequence, 610  
 nested, 385
- set, 26, 625, 631  
 active, 140, 244, 631  
 Cartesian product, 40

- closed, 32–35, 60, 61, 143
  - and open, 32, 35
- cone, 90, 96, 114, 116, 117, 129, 130, 132, 143, 229, 230, 416
  - exposed, 426
  - polyhedron, 114, 116
- connected, 31
- convex, 31
  - invariance, 39, 40
  - projection on, 135, 595, 596
- dense, 76
- difference, 40, 130, 624
- empty, 33, 34, 39, 82, 529
  - affine, 33
  - closed, 35
  - cone, 77
  - face, 75
  - hull, 51, 53, 58
  - open, 35
- feasible, 7, 39, 57, 135, 158, 171, 221, 631
  - dual, 228
  - extreme point, 248
  - primal, 224, 228
- intersection, *see* intersection
- invariance, 394
- level, 159, 183, 184, 192, 199, 208, 221
  - affine, 183
- minimal, 53, 63, 114, 320, 627
  - affine, 584
  - extreme, 87, 145
  - halfspace, 114, 115
  - hyperplane, 114
  - orthant, 170
  - positive semidefinite cone, 92, 141, 202
- nonconvex, 26, 35, 77–81, 402
- nullspace, 138
- open, 32–35, 131
  - affine, 33
  - and closed, 32, 35
- ordered, 26
- origin, 32, 77, 82, 85, *see* zero
- orthogonal, 138
- projection of, 42
- Schur-form, 93
- smooth, 388
- solution, 171, 333, 631
- sublevel, 185, 199, 204, 207
  - convex, 186, 199, 500
  - nested, 199
  - normal, 159, 193
  - quasiconvex, 207
- sum, *see* sum, vector
- superlevel, 186, 207
  - quasiconcave, 207
- union, 39, 76, 77, 79, 99
- sgn, 208, 240, 305, 509, 558, 559, 565, 626, 633
- Shannon, 299
- shape, 21
- shell, 260
- Sherman-Morrison-Woodbury, 517
- shift, *see* invariance, translation
- shroud, 372, 373, 429
  - cone, 142, 452
- SIAM, 222
- sigma delta, 19, 20
- sign, 633
- signal
  - distortion, *see* distortion
  - dropout, 276–280
  - processing, 20, 21
    - analog, 209, 215
    - compression, 330
    - digital, 19, 20, 28, 217, 276–280, 299, 400
    - reflection, 20, 364
- similarity, 400, *see* transformation
- simplex, 119, 120
  - area, 410
  - content, 411
  - method, 8, 223, 266, 282, 354
  - nonnegative, 280–282, 604
  - regular, 120, 411
  - unit, 119, 120
  - volume, 411
- singular value, 506, 530
  - $\sigma$ , 485, 486, 626
  - constraint, 58
  - decomposition, 506–510, 530
    - compact, 506
    - convex iteration, 336, 337
    - diagonal, 540
    - ellipse, 507, 508
    - full, 507
    - geometrical, 508
    - positive semidefinite, 509
    - pseudoinverse, 510
    - real, 506
    - subcompact, 506
    - symmetric, 509, 510, 539
    - unique, 336, 338
  - eigenvalue, 506, 509, 510
  - inequality, 539
  - inverse, 510
  - largest, 58, 191, 530, 604, 635
  - normal matrix, 44, 506
  - product, 509
  - smallest, 530
  - sum of, 44, 56, 530, 531
  - triangle inequality, 56
- SIOPT, 223
- slab, 33, 260, 583
- slack
  - complementary, 233
  - variable, 222, 228, 244
- Slater, 127, 167, 229, 232, 233
- slice, 103, 141, 363, 551, 552
- slope, 193
- SNR, 630
- solid, 116, 120, 409
- solution, 221
  - analytical, 189, 221, 529
  - feasible, 68, 135, 229, 608–611, 631
    - dual, 232, 233
    - strictly, 229, 232
  - global, 221
  - local, 221
  - numerical, 264
  - optimal, 68, 71–73, 171, 631
  - problem statement as, 8, 221
  - set, 171, 333, 631
  - trivial, 32

- unique, *see* unique
- vertex, 68, 177, 223, 236, 282, 290, 354
- sort
  - function, 178, 272, 462–465, 481, 486, 626
  - index matrix, 403
  - largest entries, 178
  - monotone nonnegative, 464
  - smallest entries, 178
- span, 52, 627
- sparsity, 173–179, 234, 238, 271–276, 285, 286, 297, 299, 315, 632
  - gradient, 300
  - nonnegative, 176, 275, 280–284
- spectahedron, 53
- spectral
  - cone, 396–398, 462, 626, 629
  - dual, 398, 399
  - orthant, 464
  - factorization, 266–268
  - inequality, 396
  - norm, *see* norm
  - projection, *see* projection
- sphere, *see* hypersphere
  - packing, 354
- square, *see* matrix, *see* norm
- square root, 622, 629, *see* function, fractional, *see* matrix
- Srebro, 56, 530
- Starck, 301
- steepest descent, 193, 551
- Stiemke, 134
- strain, *see* problem
- Strang, 492
- stress, *see* problem
- strict
  - complementarity, 233
  - feasibility, 229, 232
  - positivity, 383
  - triangle inequality, 385
- Sturm, 514
- subject to, 633
- submatrix
  - principal, 357, 408, 411, 425
  - face, 98
  - leading, 99, 383, 385
  - positive semidefinite, 497
  - rank, 98, 497
- subset, proper, 623
- subspace, 32, 627
  - 0, 32, 47, 627, 628, 632
  - algebra, 32, 69, 515
  - antihollow
    - antisymmetric, 48, 49, 456, 628
    - symmetric, 48, 628
  - antisymmetric, 47
  - Cartesian, 273, 603
  - complementary, *see* complement
  - fundamental, 41, 69, 516, 519, 520, 568, 575
    - projector, 575
  - geometric center, 374, 375, 415, 439, 443–445, 470, 513, 593, 628
    - dimension, 374, 593
    - orthogonal complement, 97, 371, 439, 593, 628
  - hollow, 48, 374, 445, 628
- dimension, 49, 374
- independence, 32, 624
- intersection, 73
  - hyperplane, 70
  - nullspace basis span, 70
  - parallel, 33
  - proper, 32, 48, 49, 125, 149, 153, 441, 600
  - representation, 69
  - smallest, face, 98
  - symmetric, 46
  - tangent, 96
  - translation invariant, 371, 375, 439, 593, 628
  - trivial, *see* subspace, 0
  - vectorization, 73
- successive
  - approximation, 607
  - projection, 607
- sufficient, 624
- sum, 40, 626
  - empty, 626
  - Minkowski, 40
  - of eigenvalues, 44, 485, 493, 533, 534
  - of extremes, 87
  - of functions, 171, 196, 218, 219
  - of matrices
    - eigenvalues, 497
    - nullspace, 73
    - rank, 107, 493, 494
    - of singular values, 44, 56, 530, 531
    - vector, 40, 518, 624
      - of cones, 81, 130
      - orthogonal, 130, 600, 624
      - unique, 47, 122, 518, 624
  - superset, 623
  - supremum, 51, 529, 633, *see* maximum
    - of affine functions, 185, 530, 532
    - of convex functions, 185, 191, 219
    - of quasiconvex functions, 219
    - supporting hyperplane, 65–67, 597
  - surjection, 631, *see* surjective
  - surjective, 45, 373, 375, 377, 463, 631
    - linear, 377, 378, 390, 391
  - SVD, *see* singular value decomposition
  - svec, 47, 222, 634
  - Swiss roll, 25
  - Sylvester, 494, 518
  - symmetry, 342
  - T –, 621
  - tangent
    - line, 37, 359
    - subspace, 96
  - tangential, 37, 67, 360
  - Tanner, 275
  - Tarazaga, 415, 420, 426, 476
  - taxicab, *see* distance, 1-norm
  - Taylor series, 200, 201, 473, 551–558, 565
  - tensor, 545–547
  - tesseract, 44
  - tetrahedron, 120, 281, 410
    - angle inequality, 409
    - regular, 409, 419
  - Theobald, 249, 493
  - theorem

- 0 eigenvalues, 510
- alternating projection, distance, 610
- alternative, 132–134, 139
  - EDM, 453
  - semidefinite, 230
  - weak, 134
- Barvinok, 110
- Bunt-Motzkin, 595
- Carathéodory, 131, 590
- compressed sensing, 275, 299
- cone
  - faces, 85
  - intersection, 81
  - conic coordinates, 165
  - convexity condition, 197
  - decomposition
    - dyad, 514
  - directional derivative, 552
  - discretized membership, 137
  - dual cone intersection, 163
  - duality
    - strong, 128, 233
    - weak, 127, 231
  - EDM, 390
  - eigenvalue
    - of difference, 497
    - of sum, 497
    - order, 496
    - zero, 510
  - elliptope vertices, 389
  - exposed, 88
  - extreme existence, 75
  - extremes, 87, 88, 119
  - Farkas' lemma, 131–134
    - not positive definite, 231
    - positive definite, 230
    - positive semidefinite, 229
  - fundamental
    - algebra, 503
    - algebra, linear, 69
    - convex optimization, 221
  - generalized inequality and membership, 131
  - Geršgorin discs, 105, 106
  - gradient monotonicity, 196
  - Hahn-Banach, 60, 65, 68
  - halfspaces, 61
  - Hardy-Littlewood-Pólya, 463
  - hypersphere, 389
  - inequalities, 463
  - intersection, 39
  - inverse image, 40, 43
    - closedness, 143
  - Klee, 87
  - line, 205, 218
  - linearly independent dyads, 518
  - majorization, 488
  - mean value, 555
  - Minkowski, 117
  - monotone nonnegative sort, 464
  - Motzkin, 595
    - transposition, 134
  - nonexpansivity, 602
  - pointed cones, 82
  - positive semidefinite, 491
    - convex subsets, 108
  - matrix sum, 107
  - principal submatrix, 497
  - symmetric, 498
  - projection
    - algebraic complement, 601
    - on affine, 579
    - on cone, 599
    - on convex set, 596
    - on PSD geometric intersection, 443
    - on subspace, 42
    - unique minimum-distance, 596
    - via dual cone, 601
    - via normal cone, 596
  - projector
    - rank trace, 574, 577
    - semidefinite, 498
  - proper-cone boundary, 85
  - Pythagorean, 291, 368, 605
  - range of dyad sum, 519
  - rank
    - affine dimension, 381
    - partitioned matrix, 501
    - Schur-form, 501, 502
    - trace, 574, 577
  - real eigenvector, 503
  - sparse sampling, 275, 299
  - sparsity, 176
  - Sylvester, 494, 518
  - Tour, 392
  - Weyl, 117
    - eigenvalue, 497
    - zero eigenvalues, 510
  - thin, *see* matrix
  - tight, 343, 632
  - Torgerson, 376
  - tr, 485, 493, 495, 634
  - trace, 183, 485, 530, 531, 561, 634
    - commutative, 493
    - derivative, 205, 562
    - eigenvalues, 485, 493
    - heuristic, 249, 250, 330, 471–473, 531
    - inequality, 495, 496
    - maximization, 439
    - minimization, 249, 250, 330, 357, 471–473, 500, 531
    - nonnegative, 495
    - of product, 43, 486, 487, 493, 495
      - gradient, 561
    - positive semidefinite, 495
    - product, 487, 495
    - projection matrix, 574, 577
    - rank gap, 471
    - vec, 43, 486, 487
    - zero, 510, 512
  - trajectory, 394, 434
    - sensor, 362
  - transform
    - cosine, discrete, 276, 301
    - Fourier, 266, 267
      - discrete, 45, 301, 306, 622, 630
      - discrete time, 217
      - inverse discrete, 301, 303
    - Karhunen-Loéve, 400
    - Laplace, 209–211, 267, 268
    - Legendre-Fenchel, 472

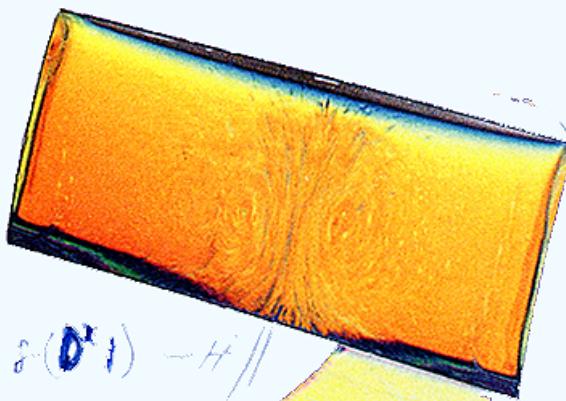
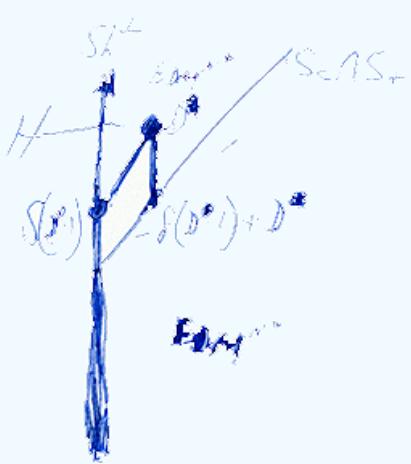
- sparsifying, 299, 301
- $z$ , 217
- transformation
  - affine, 40, 42, 65, 143
  - function, 196, 218
  - inverse, 42
  - positive semidefinite cone, 93
  - bijective, *see* bijective
  - bilinear, 217
  - congruence, 395, 494
  - coordinate, 417, 418
    - rotation, 467, 527
  - injective, *see* injective
  - invertible, 42, 44, 45
  - linear, 32, 41, 64, 112, 568
    - cone, 81, 85, 112, 116, 132, 143, 145
    - cone, dual, 130, 132, 450
    - inverse, 41, 416, 568
    - polyhedron, 112, 116, 132
  - rank-one, 334
  - rigid, 21, 360
  - similarity, 496, 590
    - sign, 498
  - surjective, *see* surjective
- transitivity
  - face, 76
  - order, 84, 91, 495, 632
- trefoil, 436, 438, 440
- triangle, 343
  - inequality, *see* inequality
- triangulation, 22, 162
- trigonometry
  - distance, 367
  - Gram, 348
  - inequality, 344
  - inner product, 42, 351
  - law of cosines, 367, 410
  - relative angle, 369
- trilateration, 22, 39
  - tandem, 363
  - unique, 357–359
- trivial, 632, *see* zero
- Tucker, 133, 161, 461
  
- $- U -$ , 53
- unbounded below, 68, 71, 72, 133, 218
- underdetermined, 332, 622
- unfolding, 436, 441
- unfurling, 436
- unimodal, 206, 207
- unique
  - cone, 90, 122, 416
  - dual, 122, 145
  - direction, extreme, 86
  - EDM, 391, 401, 462
  - eigen
    - value, 102, 503
    - vector, 504
  - expansion, biorthogonal, 147–150, 426, 518, 577, 580–582
  - generators, 119
  - Gram matrix, 371, 375
  - hull
    - affine, 50
  
- $- V -$ , 350
- value
  - absolute, 173, 201
  - eigen, *see* eigen
  - objective, 231, 241
  - saddle, 127, 129, 223
- Van Loan, 489
- variable
  - dual, 125, 161, 162, 167, 625
  - matrix, 63
  - nonnegative, *see* constraint
  - slack, 222, 228, 244
- variation
  - calculus, 135
  - total, 304
- vec, 42, 94, 302, 486, 487, 634
  - product, 487
  - trace, 43, 486, 487
- vector, 32, 631
  - binary, 51, 237, 292, 324, 326, 388, 625, 627
  - difference, 40, 130, 624
  - direction, *see* direction
  - dual, 616
  - entry, *see* entry
  - indices, 621
  - normal, 61, 595, 619
  - optimization, 83, 84, 171, 172, 270, 304
  - parallel, 578, 579, 581, 582
  - Perron, 395
  - point, 87, 631
  - primal, 616
  - product
    - inner, 42–44, 60, 153, 209, 223, 299, 367, 368, 492, 525, 587, 625
    - inner, positive semidefinite, 141
    - outer, *see* matrix, dyad
    - zero, *see* zero, product
  - pseudoinverse, 568
  - quadrature, 112, 372
  - reflection, 526

- rotation, 105, 372, 373, 525
- space, 31, 32, 627
  - ambient, 32, 441, 451
  - Euclidean, 341
  - sum, *see* sum
- vectorization, 42, *see* vec svec & dvec
  - inner product, 42–44, 492
  - Kronecker product, 487
  - projection, 430, 587, 588, 593
  - subspace, 73
  - symmetric, 47
  - hollow, 49
- Venn diagram
  - EDM, 457
  - program class, 224
  - sets, 117
- vertex, 36, 37, 57, 75, 78, 177, 388
  - description, *see* description
  - elliptope, 236, 388, 389
  - Laurent, 388
  - of cone, 82
    - none, 82, 86, 87, 114, 115, 119, 128
  - polyhedron, 57, 117, 177
  - solution, 68, 177, 223, 236, 282, 290, 354
- Vetterli, 299
- volume, *see also* content
  - facet, 410, 411
  - inequality, 411
  - polyhedron, 411
  - pyramid, 412
  - simplex, 411, 412
  - tetrahedron, 411
- von Neumann, 324, 477, 539, 607, 608, 611
- Voronoi diagram, 365
- vortex, 423
  
- $W$  —, 248
- wavelet, 300, 301
- wedge, 86, 125
- Weinberger & Saul, 25, 436, 438, 440
- Wells, 189, 415, 420, 426, 476, 477
- Weyl, 117, 132, 497
- wide, *see* matrix
- Wiener, 607
- wireless location, 21, 252, 259, 364
- womb, 544
- Woodbury, 517
- Wright, 221
- Wüthrich, 23
  
- $X$  —, 627
  
- $Y$  —, 205
- Yates, 503
- Ye, 134, 228, 361, 366
- Youla, 604
- Young, 27, 460, 462, 484, 532
  
- $Z$  —, 61
- $z$ 
  - transfer function, 217
- transform, 217
- Zenodorus, 529
- zero, 510–514, 632, *see* origin
  - boundary, 82, 85
  - cone, 82, 85
  - definite, 513, 514
  - diagonal, 510, 569
  - eigenvalue, 510, 516
  - entry, 510, 569
  - function, 171
  - gradient, 135, 193, 194, 592
  - interior, 82, 85
  - matrix, 316
  - norm, *see* 0-norm
  - norm-, 510
  - pad, 241, 507
  - product, 42, 61, 63, 73, 132, 512–514, 628
    - PSD, 91, 97–101, 233, 512, 513
  - pseudoinverse, 569
  - rank, 97, 109, 225, 226, 250
  - solution, 32, 632
  - subspace, 32, 47, 627, 628, 632
  - trace, 510, 512
  - transfer function, 210
- Zhang
  - Fuzhen, 490, 491
  - Shuzhong, 514
- Ziegler, 354
- Zinoviev, 356
- Zolotarev, 356





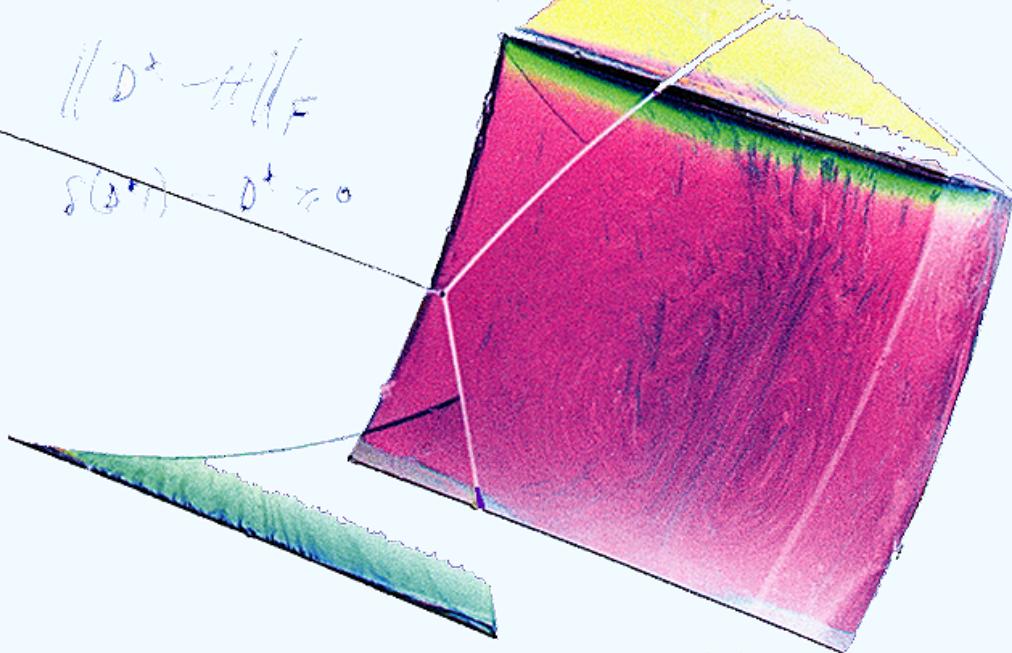
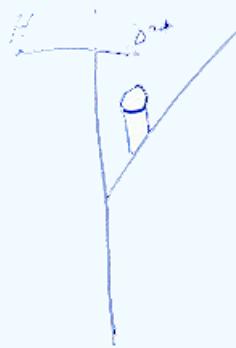




$$\| D^z - \delta(D^z) - H \|$$

$$\min \| D^z - H \|_F$$

$$\delta(D^z) = D^z \neq 0$$

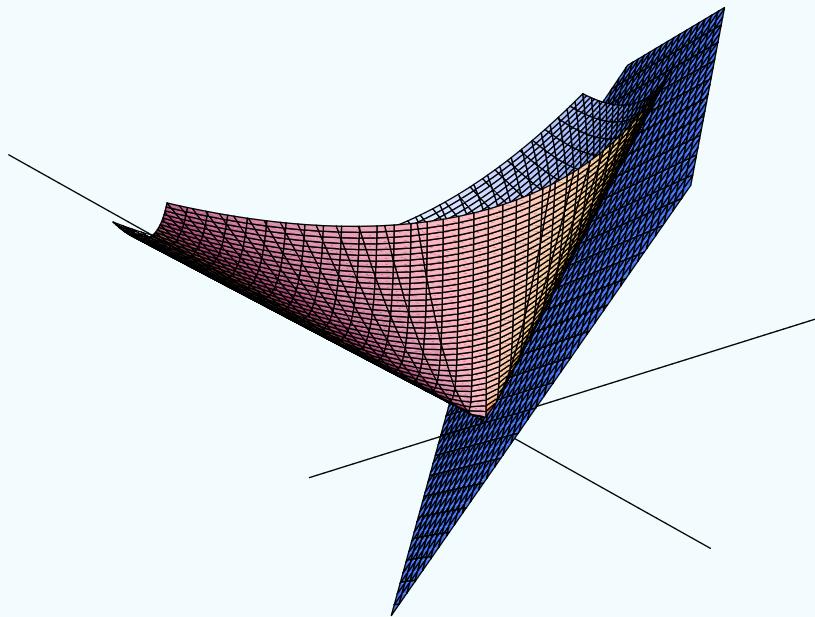


$$\delta(A) = A - \tilde{\delta}(A)$$

$$S^z S^z \delta(A) \xrightarrow{S^z S^z} f(D^z)$$

~~$$f(D^z) * f(D^z)$$~~

~~$$D^z - \delta(D^z) = (D^z -$$~~



Dattorro

*Convex Optimization † Euclidean Distance Geometry  $2\varepsilon$*

