

Bregman proximal methods for semidefinite optimization

Lieven Vandenberghe (UCLA)

joint work with Hsiao-Han Chao and Xin Jiang

Intersections between Control, Learning, and Optimization

IPAM, February 28, 2020

Semidefinite program (SDP)

$$\begin{aligned} & \text{minimize} && \text{tr}(CX) \\ & \text{subject to} && \text{tr}(A_i X) = b_i, \quad i = 1, \dots, m \\ & && X \succeq 0 \end{aligned}$$

X is a symmetric $n \times n$ matrix; $X \succeq 0$ means X is positive semidefinite

Interior-point methods

- general-purpose implementations for dense problems do not scale well
- each iteration involves computations with complexity m^3 , m^2n^2 , mn^3
- customization to exploit problem structure is difficult

Proximal splitting methods (ADMM, primal–dual hybrid gradient, ...)

- exploiting structure in linear constraints is straightforward
- require eigenvalue decompositions for projections on positive semidefinite cone

Sparse semidefinite programs

large SDPs often have sparse coefficient matrices C, A_1, \dots, A_m

- relaxations of combinatorial graph optimization problems
- semidefinite relaxations of polynomial optimization problems

Example: relaxation of maximum-cut problem

$$\begin{aligned} & \text{maximize} && \text{tr}(LX) \\ & \text{subject to} && X_{ii} = 1, \quad i = 1, \dots, n \\ & && X \geq 0 \end{aligned}$$

L is weighted graph Laplacian

- complexity of general-purpose interior-point solver: $O(n^4)$ per iteration
- customized interior-point solver: $O(n^3)$ per iteration
- proximal splitting method: $O(n^3)$ per iteration (projection on p.s.d. cone)

Nonnegative trigonometric polynomials

$$F_x(\omega) = x_0 + \sum_{k=1}^n (x_k e^{-jk\omega} + \bar{x}_k e^{jk\omega}) \geq 0 \quad \text{for all } \omega \quad (j = \sqrt{-1})$$

- coefficients x form a semidefinite-representable convex cone K
- dual cone K^* is cone of positive semidefinite Toeplitz matrices

Applications

- source of many SDP applications in signal processing since 1990s
- recent applications to superresolution, gridless compressed sensing
- SDP formulations extend to matrix polynomials, rational (Popov) functions, ...

Complexity: convex optimization over K or K^*

- general-purpose interior-point solvers: $O(n^4)$ per iteration
- customized interior-point solvers: $O(n^3)$ per iteration
- proximal splitting methods: $O(n^3)$ per iteration (for projection on p.s.d. cone)

Outline

1. Proximal methods with generalized (Bregman) distances
2. Itakura–Saito distance for nonnegative trigonometric polynomials
3. Logarithmic barrier distance for sparse p.s.d. completable matrices

Proximal mapping

Proximal mapping: for closed convex function f

$$\text{prox}_f(y) = \underset{x}{\text{argmin}} \left(f(x) + \frac{1}{2} \|x - y\|_2^2 \right)$$

Proximal algorithms

- proximal point method: $x_{k+1} = \text{prox}_{\tau f}(x_k)$
- proximal gradient method for minimizing $f(x) + g(x)$, with g differentiable:

$$\begin{aligned} x_{k+1} &= \text{prox}_{\tau f}(x_k - \tau \nabla g(x_k)) \\ &= \underset{x}{\text{argmin}} \left(f(x) + g(x_k) + \langle \nabla g(x_k), x - x_k \rangle + \frac{1}{2\tau} \|x - x_k\|_2^2 \right) \end{aligned}$$

- ADMM, Douglas–Rachford splitting, Spingarn’s method
- primal–dual hybrid gradient (Chambolle–Pock) method

Proximal algorithms with generalized distances

- use a generalized distance $d(x, y)$ instead of $\frac{1}{2}\|x - y\|_2^2$
- for example, in proximal gradient method for minimizing $f(x) + g(x)$:

$$x_{k+1} = \operatorname{argmin}_x \left(f(x) + g(x_k) + \langle \nabla g(x_k), x - x_k \rangle + \frac{1}{\tau} d(x, x_k) \right)$$

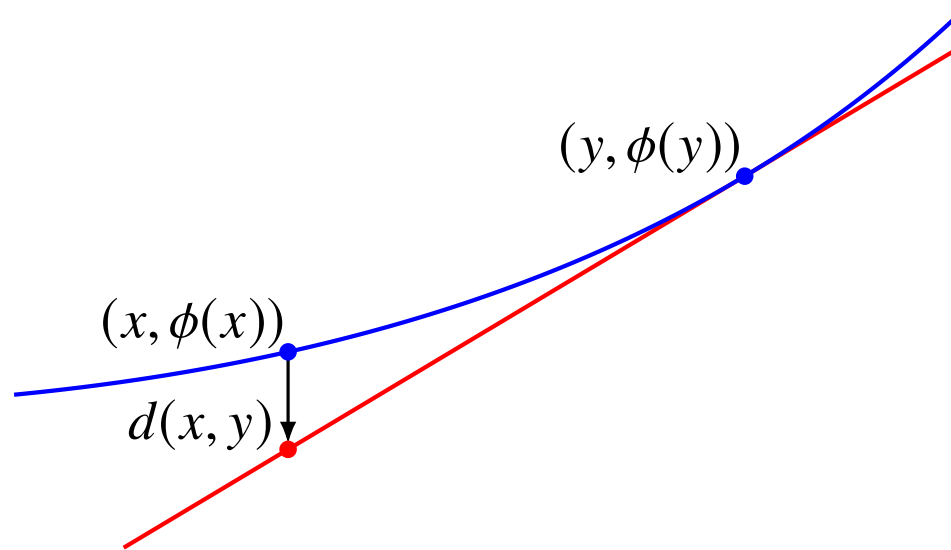
Potential benefits

1. “pre-conditioning”: use a more accurate model of $g(x)$ around x_k
2. make the generalized proximal mapping (minimizer x) easier to compute

goal of 1 is to reduce number of iterations; goal of 2 is to reduce cost per iteration

Bregman distance

$$d(x, y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle$$



- ϕ is the *kernel function*
- ϕ is convex and continuously differentiable on $\text{int}(\text{dom } \phi)$
- domain of ϕ may include its boundary or a subset of its boundary
- we define the domain of d as $\text{dom } d = \text{dom } \phi \times \text{int}(\text{dom } \phi)$

other properties of ϕ may be required (e.g., strict convexity or strong convexity)

Generalized proximal mapping

- proximal mapping of f for Bregman distance d

$$\text{prox}_f^d(y, a) = \underset{x}{\operatorname{argmin}} (f(x) + \langle a, x \rangle + d(x, y))$$

- for $d(x, y) = \frac{1}{2}\|x - y\|_2^2$, this is the standard proximal mapping

$$\begin{aligned}\text{prox}_f^d(y, a) &= \underset{x}{\operatorname{argmin}} (f(x) + \langle a, x \rangle + \frac{1}{2}\|x - y\|_2^2) \\ &= \underset{x}{\operatorname{argmin}} (f(x) + \frac{1}{2}\|x - y + a\|_2^2) \\ &= \text{prox}_f(y - a)\end{aligned}$$

Requirements

- minimizer x exists and is unique for all $y \in \text{int}(\text{dom } \phi)$ and all a
- minimizer x is in interior of $\text{dom } \phi$
- minimizer is inexpensive to compute

Example: relative entropy

$$d(x, y) = \sum_{i=1}^n (x_i \log(x_i/y_i) - x_i + y_i), \quad \text{dom } d = \mathbf{R}_+^n \times \mathbf{R}_{++}^n$$

- the Bregman distance for

$$\phi(x) = \sum_{i=1}^n x_i \log x_i, \quad \text{dom } \phi = \mathbf{R}_+^n$$

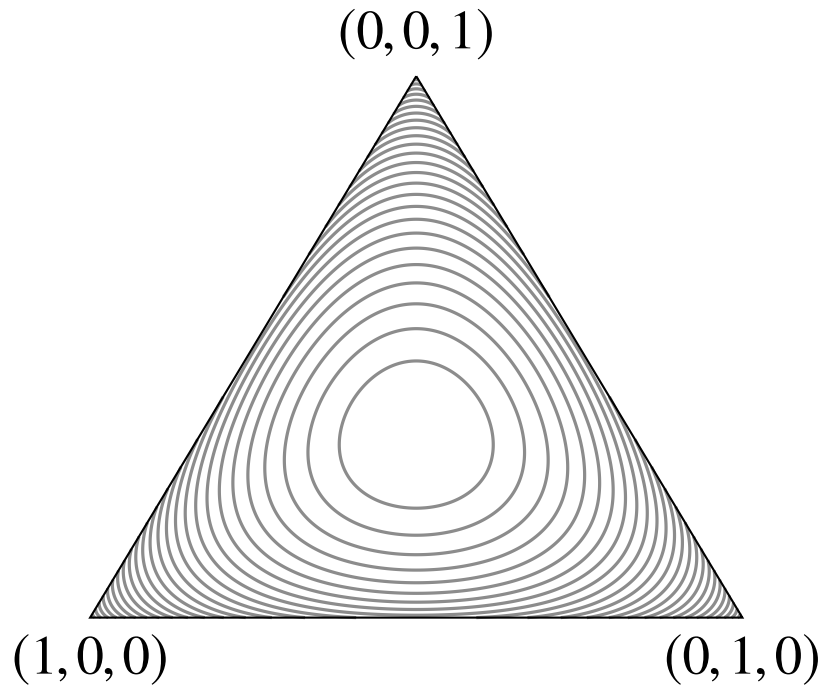
- generalized projection (proximal operator for indicator) on $H = \{x \mid \mathbf{1}^T x = 1\}$

$$\underset{\mathbf{1}^T x = 1}{\operatorname{argmin}} (a^T x + d(x, y)) = \frac{1}{\sum_{j=1}^n y_j e^{-a_j}} \begin{bmatrix} y_1 e^{-a_1} \\ \vdots \\ y_n e^{-a_n} \end{bmatrix}$$

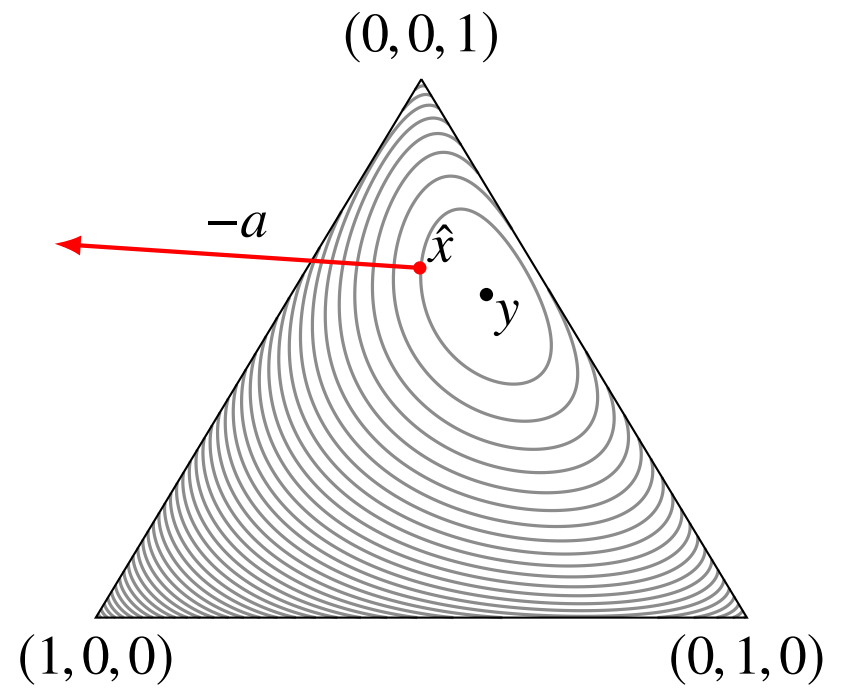
used in entropic proximal point method, exponential method of multipliers

Example: relative entropy

Contour lines of $\phi(x)$



Contour lines of $d(x, y)$



$$\hat{x} = \operatorname{argmin}_{\mathbf{1}^T x = 1} (a^T x + d(x, y))$$

Outline

1. Proximal methods with generalized (Bregman) distances
2. **Itakura–Saito distance for nonnegative trigonometric polynomials**
3. Logarithmic barrier distance for sparse p.s.d. completable matrices

Optimization over nonnegative trigonometric polynomials

we consider constraints on $x = (x_0, \dots, x_n)$:

$$x \in K, \quad x_0 = 1$$

- $K = \{x \mid F_x(\omega) \geq 0 \ \forall \omega\}$ is cone of nonnegative trigonometric polynomials
- F_x is trigonometric polynomial with coefficients x (real for simplicity)

$$F_x(\omega) = x_0 + 2x_1 \cos \omega + \dots + 2x_n \cos n\omega$$

- constraint $x_0 = 1$ normalizes $\frac{1}{2\pi} \int_0^{2\pi} F_x(\omega) d\omega = 1$

Semidefinite representation: $K = \{D(X) \mid X \geq 0\}$

$D : \mathbf{S}^{n+1} \rightarrow \mathbf{R}^{n+1}$ maps matrix to vector of diagonal sums

$$D(X)_k = \sum_{i=0}^{n-k} X_{i,i+k}, \quad k = 0, 1, \dots, n$$

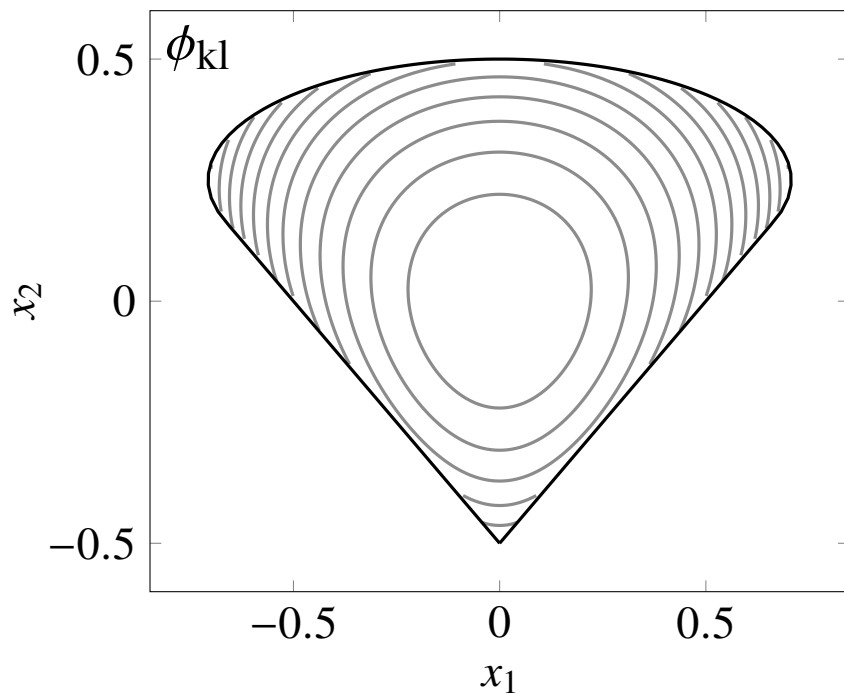
Kernel functions

kernels for Kullback–Leibler distance and Itakura–Saito distance

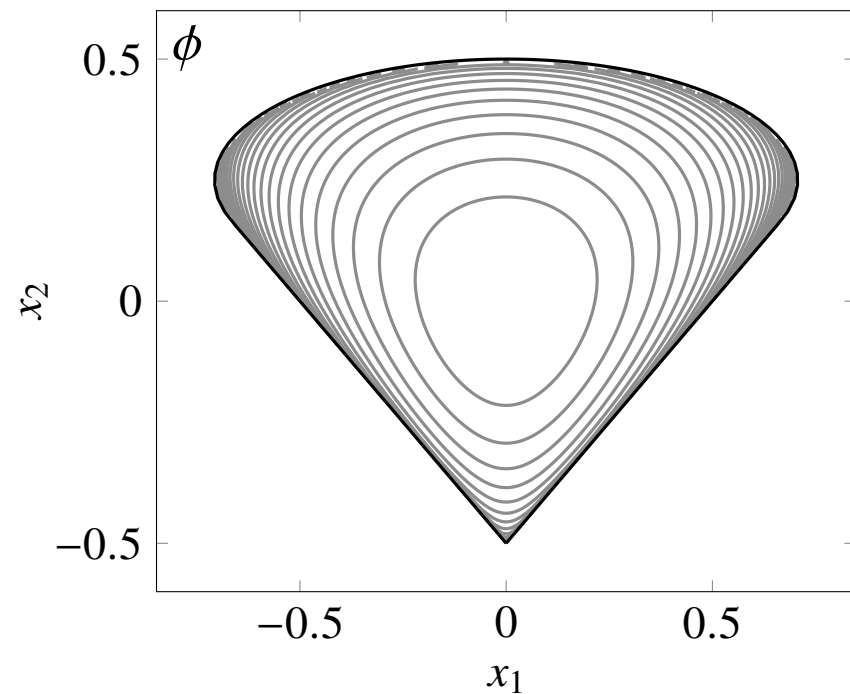
$$\phi_{\text{kl}}(x) = \frac{1}{2\pi} \int_0^{2\pi} F_x(\omega) \log F_x(\omega) d\omega$$

$$\phi(x) = -\frac{1}{2\pi} \int_0^{2\pi} \log F_x(\omega) d\omega$$

Kullback–Leibler



Itakura–Saito



- plots show contour lines on section $\{x \in K \mid x_0 = 1\}$
- ϕ is *essentially smooth*; ϕ_{kl} is not

Semidefinite representation of entropy kernel ϕ

$$\begin{array}{ll} \text{minimize (over } X) & -\log X_{00} \\ \text{subject to} & D(X) = x \\ & X \succeq 0 \end{array}$$

- for $x \in K \setminus \{0\}$, optimal value is

$$\phi(x) = -\frac{1}{2\pi} \int_0^{2\pi} \log F_x(\omega) d\omega$$

- optimal X has rank one: $X = bb^T$
- b is minimum-phase spectral factor ($b_0 + b_1z^{-1} + \dots + b_nz^{-n} \neq 0$ for $|z| > 1$)
- b is efficiently computed by spectral factorization of x : solve quadratic equation

$$D(bb^T) = x$$

Dual of semidefinite representation of ϕ

$$\text{maximize (over } y) \quad -\psi(y) - \langle x, y \rangle + 1$$

- convex function ψ is defined as

$$\psi(y) = \log(e^T T(y)^{-1} e), \quad \text{dom } \psi = \{y \mid T(y) \succ 0\}$$

where $e = (1, 0, \dots, 0)$ and $T(y)$ is the Toeplitz matrix constructed from y

- by duality, optimal value is $\phi(x)$
- optimal y is $y = -\nabla\phi(x)$, satisfies

$$T(y)b = e$$

y can be computed from spectral factor b by reverse Levinson algorithm

Itakura–Saito distance and projection

$$d(x, y) = \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{F_x(\omega)}{F_y(\omega)} - \log \frac{F_x(\omega)}{F_y(\omega)} - 1 \right) d\omega$$

- proposed in 1970s as spectral distance measure in speech processing
- generalized projection on hyperplane $H = \{x \mid x_0 = 1\}$:

$$\begin{aligned} \text{prox}_{\delta_H}^d(y, a) &= \underset{x_0=1}{\text{argmin}} (\langle a, x \rangle + d(x, y)) \\ &= \underset{x_0=1}{\text{argmin}} (\langle c, x \rangle + \phi(x)) \quad (\text{where } c = a - \nabla \phi(y)) \end{aligned}$$

- dual problem (scalar variable λ is multiplier for constraint $x_0 = 1$)

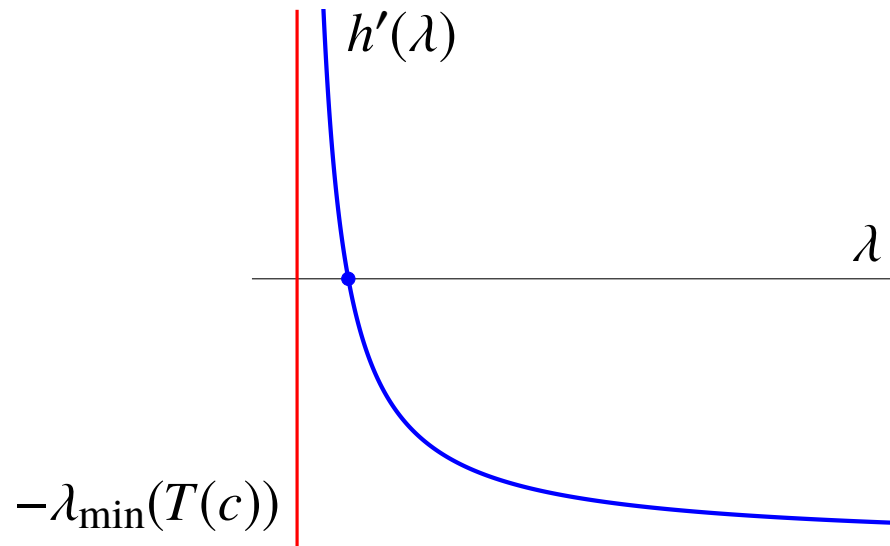
$$\text{maximize} \quad -\log \left(e^T (T(c) + \lambda I)^{-1} e \right) - \lambda$$

$e^T (T(c) + \lambda I)^{-1} e$ is the 1st element of the inverse of Toeplitz matrix $T(c) + \lambda I$

Computing Itakura–Saito projection

solve dual problem for λ , for example, by Newton's method

$$\text{maximize } h(\lambda) = -\log \left(e^T (T(c) + \lambda I)^{-1} e \right) - \lambda$$

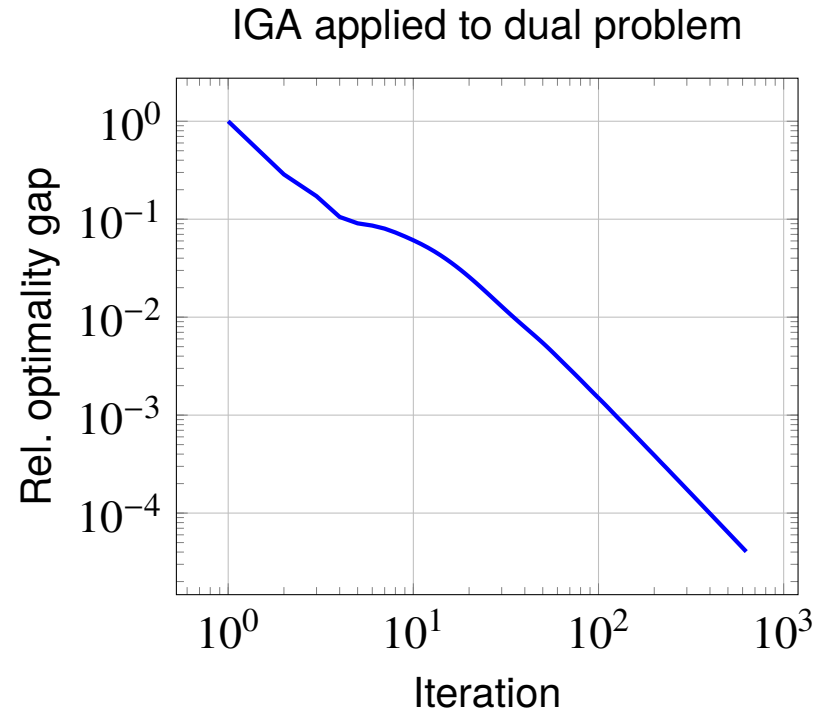
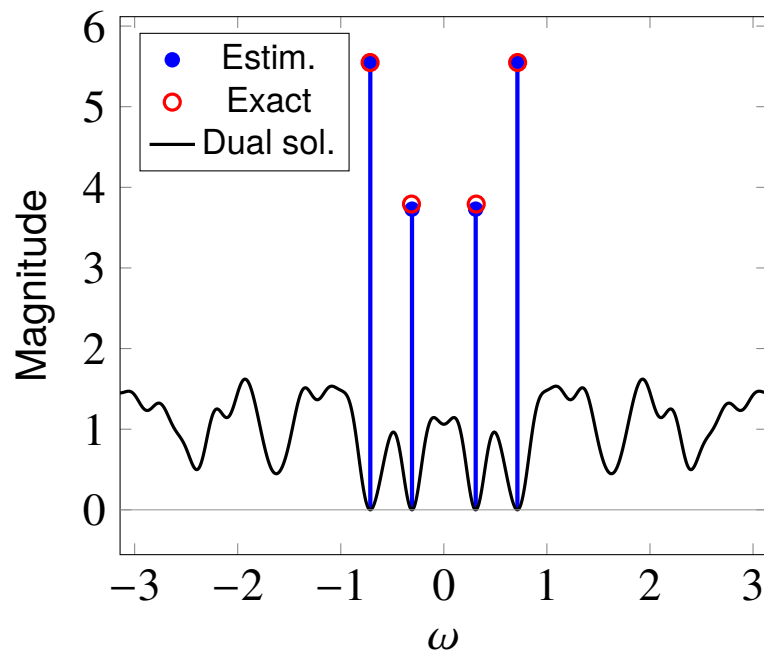


- at each Newton step, factorize positive definite Toeplitz matrix $T(c) + \lambda I$
- complexity: $O(n^2)$ with Levinson algorithm, $O(n(\log n)^2)$ with superfast solvers
- from optimal λ , compute solution $x = (1/b_0)D(bb^T)$ where $b = (T(c) + \lambda I)^{-1}e$

Covariance estimation

$$\begin{aligned} & \text{minimize (over } y, s) && \|T(y) + sI - R\|_F^2 + \gamma \text{tr}(T(y)) \\ & \text{subject to} && T(y) \succeq 0 \end{aligned}$$

- estimate parameters in signal model $v(t) = \sum_{k=1}^{\rho} c_k e^{j\omega_k t} + \text{white noise}$
- fit covariance $T(y) + sI$: low-rank p.s.d. Toeplitz plus multiple of identity
- R is sample covariance matrix ($n + 1 = 30$)



IGA: proximal gradient algorithm with Bregman distances

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & x \in C \end{array}$$

f convex with Lipschitz continuous gradient, C a convex set

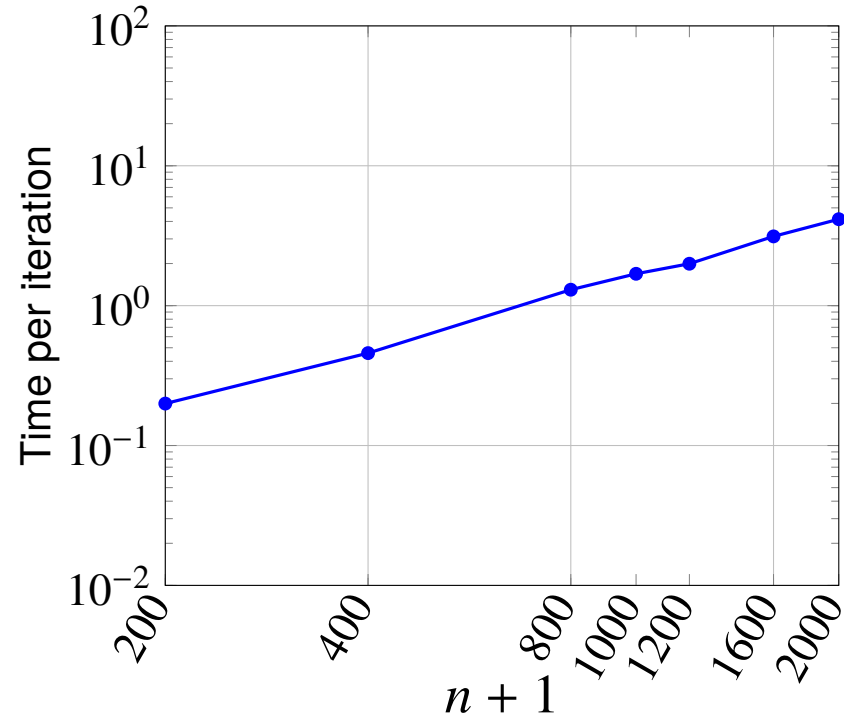
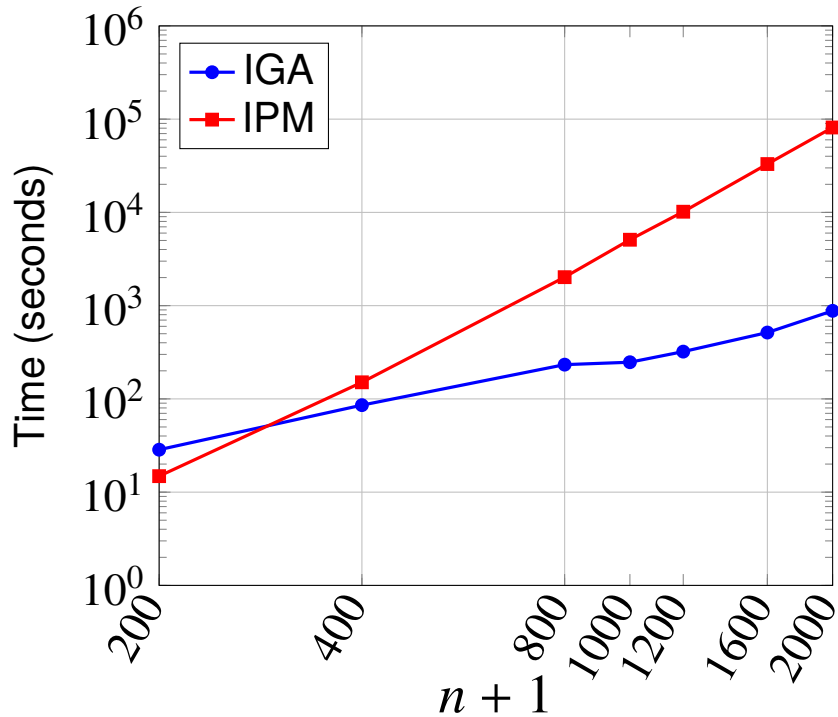
Improved gradient algorithm (IGA) [Auslender and Teboulle 2006]

$$\begin{aligned} y_{k+1} &= (1 - \theta_k)x_k + \theta_k v_k \\ v_{k+1} &= \operatorname{argmin}_{x \in C} (\langle \nabla f(y_{k+1}), x \rangle + \frac{1}{\tau_k} d(x, v_k)) \\ x_{k+1} &= (1 - \theta_k)x_k + \theta_k v_{k+1} \end{aligned}$$

extends Nesterov's 1988 accelerated projected gradient algorithm

Euclidean projection

$$\begin{aligned} &\text{minimize} && \sum_{k=0}^n (x_k - a_k)^2 \\ &\text{subject to} && x \in K, \quad x_0 = 1 \end{aligned}$$



- IPM is SDPT3/SeDuMi via CVX; IGA is Auslender–Teboulle algorithm
- number of IGA iterations is 100–200 to reach relative accuracy 10^{-4}
- about 10 Newton steps per projection; Toeplitz solver is Levinson algorithm

Outline

1. Proximal methods with generalized (Bregman) distances
2. Itakura–Saito distance for nonnegative trigonometric polynomials
3. **Logarithmic barrier distance for sparse p.s.d. completable matrices**

Sparse semidefinite program

$$\begin{array}{ll} \text{minimize} & \text{tr}(CX) \\ \text{subject to} & \text{tr}(A_i X) = b_i, \quad i = 1, \dots, m \\ & X \geq 0 \end{array}$$

- C, A_1, \dots, A_m are sparse with common sparsity pattern E
- without loss of generality, we assume E is chordal (a filled Cholesky pattern)
- optimal X is typically dense, even for sparse coefficients C, A_1, \dots, A_m

Equivalent conic linear program

$$\begin{array}{ll} \text{minimize} & \text{tr}(CX) \\ \text{subject to} & \text{tr}(A_i X) = b_i, \quad i = 1, \dots, m \\ & X \in K \end{array}$$

- variable X is a *sparse* matrix with sparsity pattern E (notation: \mathbf{S}_E^n)
- K is cone of matrices in \mathbf{S}_E^n that have a positive semidefinite completion

Centering problem

Logarithmic barrier

$$\phi(X) = \sup_{S \in \text{int } K^*} (-\text{tr}(XS) + \log \det S)$$

- dual cone K^* is cone of positive semidefinite matrices in \mathbf{S}_E^n
- ϕ is conjugate barrier of log-det barrier $\phi_*(S) = -\log \det S$ for K^*

Centering problem

$$\begin{aligned} & \text{minimize} && \text{tr}(CX) + \mu\phi(X) \\ & \text{subject to} && \text{tr}(A_i X) = b_i, \quad i = 1, \dots, m \end{aligned}$$

- solutions for $\mu > 0$ form the central path of the SDP
- optimal X is (μn) -suboptimal for the SDP

Distance generated by barrier kernel

$$\phi(X) = \sup_{S \in \text{int } K^*} (-\text{tr}(XS) + \log \det S)$$

- optimal \hat{S}_X is (sparse) inverse of max. determinant pos. definite completion of X

$$\phi(X) = \log \det \hat{S}_X - n$$

- gradient $\nabla \phi(X) = -\hat{S}_X$
- for chordal E : efficient algorithms for computing \hat{S}_X given X
- complexity is comparable with sparse Cholesky factorization with pattern E

Distance

$$\begin{aligned} d(X, Y) &= \phi(X) - \phi(Y) - \text{tr}(\nabla \phi(Y)(X - Y)) \\ &= -\log \det(\hat{S}_Y \hat{S}_X^{-1}) + \text{tr}(\hat{S}_Y \hat{S}_X^{-1}) + n \end{aligned}$$

the relative entropy (Kullback–Leibler divergence) between \hat{S}_Y and \hat{S}_X

Bregman proximal operator for centering problem

$$\begin{aligned} & \text{minimize} && \text{tr}(CX) + \mu\phi(X) \\ & \text{subject to} && \text{tr}(A_i X) = b_i, \quad i = 1, \dots, m \\ & && \text{tr} X = 1 \end{aligned}$$

- centering objective, restricted to $\text{tr} X = 1$ (alternatively, $\text{tr} X \leq 1$)

$$f(X) = \text{tr}(CX) + \mu\phi(X) + \delta_H(X), \quad H = \{X \mid \text{tr} X = 1\}$$

- Bregman proximal operator for centering objective

$$\begin{aligned} \hat{X} &= \underset{X}{\operatorname{argmin}} (f(X) + \text{tr}(DX) + \frac{1}{\tau}d(X, Y)) \\ &= \underset{\text{tr} X=1}{\operatorname{argmin}} (\text{tr}(BX) + \phi(X)) \quad \text{where } B = \frac{1}{1 + \mu\tau}(\tau(D + C) + \hat{S}_Y) \in \mathbf{S}_E^n \end{aligned}$$

- dual problem (scalar variable λ is multiplier for $\text{tr} X = 1$):

$$\text{maximize} \quad \log \det(B + \lambda I) - \lambda$$

Algorithm for Bregman proximal operator

$$\begin{array}{ll} \text{minimize} & \text{tr}(BX) + \phi(X) \\ \text{subject to} & \text{tr } X = 1 \end{array}$$

- use Newton's method to find unique solution λ of the nonlinear equation

$$\text{tr}((B + \lambda I)^{-1}) = 1 \quad (\text{with } B + \lambda I > 0)$$

- from λ , compute solution \hat{X} as projection $\Pi_E((B + \lambda I)^{-1})$ on \mathbf{S}_E^n
- for chordal sparsity patterns E , efficient algorithms exist for computing

$$g(\lambda) = \text{tr}((B + \lambda I)^{-1}), \quad g'(\lambda) = -\text{tr}((B + \lambda I)^{-2}), \quad \hat{X} = \Pi_E((B + \lambda I)^{-1})$$

from sparse Cholesky factorization of $B + \lambda I$

complexity \approx # Newton iterations \times cost of sparse Cholesky factorization

Maximum-cut problem

$$\begin{aligned} & \text{maximize} && \text{tr}(LX) \\ & \text{subject to} && \text{diag}(X) = \mathbf{1}, X \geq 0 \end{aligned}$$

- compute approximate solution on central path (parameter $\mu = 0.001/n$)
- Bregman variant of primal–dual hybrid gradient algorithm [Chambolle & Pock 2016]
- four problems from SDPLIB, four graphs from SuiteSparse collection

	n	time per prox-evaluation	Newton iterations per prox-evaluation	PDHG iterations
maxG51	1000	0.13	2.45	267
maxG32	2000	0.21	1.56	240
maxG55	5000	0.65	2.10	249
maxG60	7000	1.57	2.55	279
barth4	6019	1.62	3.57	346
tuma2	12992	2.14	4.36	375
biplane-9	21701	2.58	2.50	287
c-67	57975	3.58	4.63	378

SDP relaxation of graph partitioning

$$\begin{aligned} & \text{minimize} && \text{tr}(P^T L P X) \\ & \text{subject to} && \text{diag}(P X P^T) = \mathbf{1}, \quad X \geq 0 \end{aligned}$$

- columns of P are sparse basis of $\{x \mid \mathbf{1}^T x = 0\}$
- Bregman PDHG for centering problem (centering parameter $\mu = 0.001/n$)
- four problems from SDPLIB, four graphs from SuiteSparse

	n	time per prox-evaluation	Newton iterations per prox-evaluation	PDHG iterations
gpp100	100	0.02	1.45	305
gpp124-1	124	0.02	1.99	392
gpp250-1	250	0.03	2.78	365
gpp500-1	500	0.07	2.99	394
delaunay_n10	1024	1.76	4.32	403
delaunay_n11	2048	2.54	4.68	420
delaunay_n12	4096	3.05	4.47	367
delaunay_n13	8192	4.98	4.39	375

Primal–dual hybrid gradient (PDHG) method

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & Ax = b \end{array}$$

f is a closed convex function

Algorithm

$$y_{k+1} = z_k + \theta_k(z_k - z_{k-1})$$

$$x_{k+1} = \underset{x}{\operatorname{argmin}} (f(x) + y_{k+1}^T Ax + \frac{1}{\tau_k} d(x, x_k))$$

$$z_{k+1} = z_k + \sigma_k(Ax_{k+1} - b)$$

- parameters θ_k , σ_k , τ_k are fixed or determined by line search
- Bregman variant of primal–dual hybrid gradient (Chambolle–Pock) method [Chambolle & Pock 2016]

Summary

Bregman proximal methods for two classes of SDP-representable constraints

Nonnegative trigonometric polynomials

- Itakura–Saito distance
- cost of generalized projection is roughly $O(n^2)$

Positive semidefinite completable sparse matrices

- distance generated by logarithmic barrier
- prox-operator for centering objective
- cost roughly on the order of cost of sparse Cholesky factorization

References

1. H.-H. Chao, L. Vandenberghe, IEEE Trans. Signal Processing, 2018.
2. X. Jiang, L. Vandenberghe, in preparation.