

# Polyhedral Analysis of Quadratic Optimization Problems with Stieltjes Matrices and Indicators

Peijing Liu <sup>\*</sup>    Alper Atamtürk <sup>†</sup>    Andrés Gómez <sup>‡</sup>  
Simge Küçükyavuz <sup>§</sup>

April 5, 2024

## Abstract

In this paper, we consider convex quadratic optimization problems with indicators on the continuous variables. In particular, we assume that the Hessian of the quadratic term is a Stieltjes matrix, which naturally appears in sparse graphical inference problems and others. We describe an explicit convex formulation for the problem by studying the Stieltjes polyhedron arising as part of an extended formulation and exploiting the supermodularity of a set function defined on its extreme points. Our computational results confirm that the proposed convex relaxation provides an exact optimal solution and may be an effective alternative, especially for instances with large integrality gaps that are challenging with the standard approaches.

## 1 Introduction

Given vectors  $\mathbf{a}, \mathbf{c} \in \mathbb{R}^n$ , a *Stieltjes* matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  (that is,  $\mathbf{Q} \succ 0$  and  $Q_{ij} \leq 0$  for  $i \neq j$ ), and a convex set  $C \subseteq \mathbb{R}^{2n}$ , consider the mixed-integer quadratic optimization (MIQO) problem

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n} \mathbf{a}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{z} + \mathbf{x}^\top \mathbf{Q} \mathbf{x} \quad (1a)$$

$$\text{s.t. } \mathbf{x} \circ (\mathbf{e} - \mathbf{z}) = \mathbf{0} \quad (1b)$$

$$(\mathbf{x}, \mathbf{z}) \in C, \quad (1c)$$

---

<sup>\*</sup>Department of Industrial and System Engineering, University of Southern California, [peijingl@usc.edu](mailto:peijingl@usc.edu).

<sup>†</sup>Department of Industrial Engineering and Operations Research, University of California, Berkeley, [atamturk@berkeley.edu](mailto:atamturk@berkeley.edu).

<sup>‡</sup>Department of Industrial and System Engineering, University of Southern California, [gomezand@usc.edu](mailto:gomezand@usc.edu).

<sup>§</sup>Department of Industrial Engineering and Management Sciences, Northwestern University, [simge@northwestern.edu](mailto:simge@northwestern.edu).

where  $\mathbf{e}$  is an  $n$ -dimensional vector of ones. Problem (1) with a Stieltjes matrix  $\mathbf{Q}$  arises naturally in statistical problems with graphical models, which we discuss in §2, and others.

A fundamental step towards solving (1) effectively is designing strong convex relaxations of the optimization problem. Several approaches have been proposed for MIQO in the literature, by exploiting low-dimensional cases [13, 16] or low-rank cases [4, 5, 15, 24, 25, 26]. In this paper, we turn our attention to a critical substructure given by

$$X_Q \stackrel{\text{def}}{=} \{(t, \mathbf{x}, \mathbf{z}) \in \mathbb{R}^{n+1} \times \{0, 1\}^n : t \geq \mathbf{x}^\top \mathbf{Q} \mathbf{x}, (1b)\}.$$

Set  $X_Q$  is the mixed-integer epigraph of a quadratic function with a Stieltjes matrix and indicators. Atamtürk and Gómez [3] show that if  $\mathbf{a}$  has all entries of the same sign, problem (1) with a Stieltjes matrix  $\mathbf{Q}$  is polynomial-time solvable. However, a tractable convex relaxation of (1) that guarantees optimality under the same conditions has been missing. This paper presents such a convex relaxation.

## Outline

The rest of the paper is organized as follows. In §2, we discuss the applications of quadratic optimization with a Stieltjes matrix and indicators. In §3, we provide the necessary background for the paper. In §4, we convexify  $X_Q$  by exploiting an underlying polyhedral structure. In §5, we present experimental results illustrating the computational impact of the proposed convexification.

## Notation

Given any  $n \in \mathbb{Z}_+$ , let  $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$ . We denote vectors and matrices in **bold**. Moreover,  $\mathbf{0}$  denotes either a vector or matrix of zeros (whose dimension is clear from the context),  $\mathbf{e}$  and  $\mathbf{E}$  denote a vector and matrix of ones, respectively. Given a set  $S \subseteq [n]$  we let  $\mathbf{e}_S \in \mathbb{R}^n$  denote the indicator vector of  $S$ , i.e.,  $(\mathbf{e}_S)_i = 1$  if  $i \in S$  and  $(\mathbf{e}_S)_i = 0$  otherwise. Moreover, given  $i \in S$ , we let  $\mathbf{e}_i \stackrel{\text{def}}{=} \mathbf{e}_{\{i\}}$  denote the  $i$ -th standard vector of  $\mathbb{R}^n$ . Similarly, we let  $\mathbf{E}_{ij}$  be the square matrix whose dimensions can be inferred from the context, with a 1 in the  $(i, j)$ -th position and 0 elsewhere.

We let  $\mathbb{S}_+^n$  denote the cone of positive semidefinite  $n \times n$  matrices, and  $\mathbb{S}_{++}^n$  the set of positive definite  $n \times n$  matrices. Given a matrix  $\mathbf{W}$ , we let  $\mathbf{W}^\dagger$  denote its pseudoinverse. A special case of the pseudoinverse that is used throughout the paper pertains to matrices of the form  $\mathbf{W} = \begin{pmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$  where  $\mathbf{A}$  is invertible, in which case  $\mathbf{W}^\dagger = \begin{pmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$ . Given two matrices  $\mathbf{W}_1$  and  $\mathbf{W}_2$  of the same dimension, we let  $\mathbf{W}_1 \circ \mathbf{W}_2$  denote their Hadamard (entrywise) product.

Given a matrix  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and set  $S$ , denote by  $\mathbf{W}_S \in \mathbb{R}^{S \times S}$  the submatrix of  $\mathbf{W}$  induced by  $S$ . Moreover, given  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $S \subseteq [n]$ , observe that  $(\mathbf{W} \circ \mathbf{e}_S \mathbf{e}_S^\top) \in \mathbb{R}^{n \times n}$  is the matrix that coincides with  $\mathbf{W}_S$  in the entries indexed by  $S$  and is 0 elsewhere. To simplify the notation, given  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $S \subseteq [n]$ , we let

$$\mathbf{W}_S^* \stackrel{\text{def}}{=} (\mathbf{W} \circ \mathbf{e}_S \mathbf{e}_S^\top)^\dagger;$$

if  $S$  corresponds to the first indices of  $[n]$ , then note that  $\mathbf{W}_S^* = \begin{pmatrix} \mathbf{W}_S^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$ .

## 2 MIQOs with Stieltjes matrices

Problem (1) with a Stieltjes matrix  $\mathbf{Q}$  arises naturally in statistical problems with graphical models, which we discuss in §2.1. Set  $X_Q$  is critical to such problems because a convex relaxation of (1) can be obtained as

$$\begin{aligned} \min_{(t, \mathbf{x}, \mathbf{z}) \in \mathbb{R}^{2n+1}} \quad & \mathbf{a}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{z} + t \\ \text{s.t.} \quad & (t, \mathbf{x}, \mathbf{z}) \in \text{conv}(X_Q), (\mathbf{x}, \mathbf{z}) \in C, \end{aligned}$$

and the relaxation is exact if  $C = \mathbb{R}^{2n}$ .

Understanding  $\text{conv}(X_Q)$  for Stieltjes matrices can also be helpful in solving MIQO problems with non-Stieltjes quadratic matrices. In such cases, a standard approach in the literature is to decompose  $\mathbf{Q}$  into simpler matrices of the form  $\mathbf{Q} = \mathbf{Q}_0 + \sum_{k \in K} \mathbf{Q}_k$  for some index set  $K$ , where  $\mathbf{Q}_k$  are “simple” matrices and  $\mathbf{Q}_0 \succeq 0$  is a remainder “complicated” matrix. Then relaxations of problem (1) can be obtained as

$$\min_{(t, \mathbf{x}, \mathbf{z}) \in \mathbb{R}^{|K|+2n}} \quad \mathbf{a}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{z} + \mathbf{x}^\top \mathbf{Q}_0 \mathbf{x} + \sum_{k \in K} t_k \quad (2a)$$

$$\text{s.t.} \quad (t_k, \mathbf{x}, \mathbf{z}) \in \text{conv}(X_{Q_k}) \quad \forall k \in K \quad (2b)$$

$$(\mathbf{x}, \mathbf{z}) \in C. \quad (2c)$$

The most prevalent relaxation in the literature to tackle (1) is the *perspective relaxation* [1, 11, 12, 14], which is a special case of the approach outlined where  $|K| = 1$  and  $\mathbf{Q}_1$  is a diagonal positive definite matrix. A good understanding of relaxations of  $X_Q$  allows one to extend the perspective relaxation to more general (non-diagonal) classes of Stieltjes matrices. In fact, as we discuss in §2.2, matrices  $\mathbf{Q}_k$  are not required to be Stieltjes to utilize the results of this paper.

### 2.1 A direct application

Mixed-integer convex quadratic problems with Stieltjes matrices arise, for example, in inference problems with Besag-York-Mollié graphical models [9]. In

this class of graphical models, the vertex set  $V$  of graph  $\mathcal{G} = (V, E)$  represents latent random variables,  $Y$ , and the edge set  $E$  represents relationships between random variables. The existence of an edge  $[i, j] \in E$  indicates that the random variables  $i$  and  $j$  should have similar values. Thus, the values of the unobserved random variables can be estimated as the optimal solution of the optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^V, \mathbf{z} \in \{0,1\}^V} \sum_{i \in V} a_i (y_i - x_i)^2 + \sum_{[i,j] \in E} a_{ij} (x_i - x_j)^2 \quad (3a)$$

$$\text{s.t. (1b) - (1c),} \quad (3b)$$

where  $y_i$  represents some noisy measurement of the value of random variable  $Y_i$ ,  $i \in V$ ,  $\mathbf{a} \geq \mathbf{0}$  are parameters to be tuned, and the constraints incorporate logical constraints on the estimators. We refer the reader to Han et al. [17] for additional information on this class of combinatorial inference problems. As mentioned in the introduction, if  $C = \mathbb{R}^n \times \{0, 1\}^n$  or  $C = \mathbb{R}_+^n \times \{0, 1\}^n$ , optimization problems with Stieltjes matrices and indicators can be solved in polynomial time. We refer the reader to Atamtürk and Gómez [3] for the case with  $\mathbf{a} \leq \mathbf{0}$  and  $C = \mathbb{R}_+ \times \{0, 1\}^n$ .

## 2.2 Stieltjes-equivalent classes

In several cases, quadratic functions with non-Stieltjes matrices can be transformed into equivalent expressions with Stieltjes matrices. Specifically, given a vector  $\mathbf{x} \in \mathbb{R}^n$ , matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  and index set  $I \subseteq \{1, \dots, n\}$ , define

$$\bar{\mathbf{x}} = \begin{cases} -x_i & \text{if } i \in I \\ x_i & \text{otherwise,} \end{cases} \quad \text{and } \bar{\mathbf{Q}}_{ij} = \begin{cases} -Q_{ij} & \text{if } |\{i, j\} \cap I| = 1 \\ Q_{ij} & \text{otherwise.} \end{cases}$$

Then observe that  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \bar{\mathbf{x}}^\top \bar{\mathbf{Q}} \bar{\mathbf{x}}$ . Thus, since  $(t, \mathbf{x}, \mathbf{z}) \in X_{\mathbf{Q}} \Leftrightarrow (t, \bar{\mathbf{x}}, \mathbf{z}) \in X_{\bar{\mathbf{Q}}}$ , we find that we can convexify sets with non-Stieltjes matrices  $\mathbf{Q}$  provided that the transformed matrix  $\bar{\mathbf{Q}}$  is Stieltjes.

Improved formulations for convexifications using Stieltjes-equivalent classes have been presented in the literature. First, there have been notable efforts in studying  $X_{\mathbf{Q}}$  where  $n = 2$  [8, 13, 18, 27]. Naturally, any positive definite  $2 \times 2$  matrix is Stieltjes-equivalent: if  $Q_{12} \leq 0$ , then  $\mathbf{Q}$  is already Stieltjes; otherwise, after the substitution  $\bar{x}_1 = -x_1$ , we recover an equivalent expression with a Stieltjes matrix. In another line of work, Liu et al. [20] describe the closure of the convex hull of  $X_{\mathbf{Q}}$  (in a polynomially-sized extended formulation) when  $\mathbf{Q}$  is tridiagonal. It can also be verified that tridiagonal matrices are Stieltjes equivalent.

### 3 Preliminaries

In this section, we discuss some preliminary results concerning the convexification of  $X_Q$  relevant to the current paper.

**Theorem 1** (Wei et al. [27]). *Given any matrix  $\mathbf{Q} \in \mathbb{S}_{++}^n$ , the closure of the convex hull of  $X_Q$  can be described in an extended formulation as*

$$\text{cl conv}(X_Q) = \left\{ (t, \mathbf{x}, \mathbf{z}) \in \mathbb{R}^{2n+1} : \exists \mathbf{W} \in \mathbb{R}^{n \times n} \text{ such that } \begin{pmatrix} \mathbf{W} & \mathbf{x} \\ \mathbf{x}^\top & t \end{pmatrix} \in \mathbb{S}_+^{n+1}, \right. \\ \left. (\mathbf{z}, \mathbf{W}) \in \text{conv}(P_Q) \right\},$$

where  $P_Q \stackrel{\text{def}}{=} \left\{ (\mathbf{z}, \mathbf{W}) \in \{0, 1\}^n \times \mathbb{R}^{n \times n} : \mathbf{W} = (\mathbf{Q} \circ \mathbf{z}\mathbf{z}^\top)^\dagger \right\}$ .

Note that  $P_Q$  is a finite set and, therefore,  $\text{conv}(P_Q)$  is a polytope. Consequently, we see from Theorem 1 that convexification of  $X_Q$  in a higher dimension reduces to the convexification of a polyhedral set, allowing for the use of theory and techniques from polyhedral theory. However, to utilize Theorem 1 effectively, a major challenge must be overcome: characterizing or approximating set  $\text{conv}(P_Q)$ , which has not been studied in the literature. Our goal in this paper is to close this gap for the special case of Stieltjes polytopes, as defined next.

**Definition 1** (Stieltjes polytope). Given a Stieltjes matrix  $\mathbf{Q}$ , the Stieltjes polytope associated with  $\mathbf{Q}$  is defined as

$$Z_Q = \text{conv} \left( \{(\mathbf{e}_S, \mathbf{Q}_S^*)\}_{S \subseteq [n]} \right). \quad (4)$$

*Example 1.* Let  $n = 3$ , define  $\mathbf{D} = \begin{pmatrix} 3 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 3 \end{pmatrix}$ ,  $\mathbf{q} = \mathbf{e}$  and  $\mathbf{Q} = \mathbf{D} - \mathbf{e}\mathbf{e}^\top = \begin{pmatrix} 2 & -1 & -1 \\ -1 & 3 & -1 \\ -1 & -1 & 2 \end{pmatrix}$ . For the Stieltjes matrix  $\mathbf{Q}$  the Stieltjes polytope  $Z_Q$  is the convex hull of the following eight points in  $\mathbb{R}^{12}$ :

$$\left( \mathbf{0}, \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \right), \left( \mathbf{e}_1, \begin{pmatrix} 1/2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \right), \left( \mathbf{e}_2, \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1/3 & 0 \\ 0 & 0 & 0 \end{pmatrix} \right), \\ \left( \mathbf{e}_3, \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1/2 \end{pmatrix} \right), \left( \mathbf{e}_{\{1,2\}}, \begin{pmatrix} 3/5 & 1/5 & 0 \\ 1/5 & 2/5 & 0 \\ 0 & 0 & 0 \end{pmatrix} \right), \left( \mathbf{e}_{\{2,3\}}, \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2/5 & 1/5 \\ 0 & 1/5 & 3/5 \end{pmatrix} \right), \\ \left( \mathbf{e}_{\{1,3\}}, \begin{pmatrix} 2/3 & 0 & 1/3 \\ 0 & 0 & 0 \\ 1/3 & 0 & 2/3 \end{pmatrix} \right), \left( \mathbf{e}_{\{1,2,3\}}, \begin{pmatrix} 5/3 & 1 & 4/3 \\ 1 & 1 & 1 \\ 4/3 & 1 & 5/3 \end{pmatrix} \right).$$

■

Since describing  $Z_Q$  requires finding a polyhedral description of the inverses of submatrices of  $Q$ , we review two technical lemmas concerning matrix inversion that will be used throughout the paper.

**Lemma 1** (Blockwise inversion, Lu and Shiou [22]). *The inverse of a non-singular square matrix  $R = \begin{pmatrix} A & B \\ C & D \end{pmatrix}$  is given by*

$$R^{-1} = \begin{pmatrix} A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1} & -A^{-1}B(D - CA^{-1}B)^{-1} \\ (D - CA^{-1}B)^{-1}CA^{-1} & (D - CA^{-1}B)^{-1} \end{pmatrix}.$$

We apply Lemma 1 to a special class of matrices. We summarize the relevant results in the following corollary.

**Corollary 1.** *The inverse of a positive definite matrix  $R = \begin{pmatrix} A & v \\ v^\top & d \end{pmatrix}$ , where  $A \in \mathbb{R}^{n \times n}$ ,  $v \in \mathbb{R}^n$  and  $d \in \mathbb{R}_+$  is given by*

$$R^{-1} = \begin{pmatrix} A^{-1} + A^{-1}v(d - v^\top A^{-1}v)^{-1}v^\top A^{-1} & -A^{-1}v(d - v^\top A^{-1}v)^{-1} \\ (d - v^\top A^{-1}v)^{-1}v^\top A^{-1} & (d - v^\top A^{-1}v)^{-1} \end{pmatrix}.$$

Moreover, letting  $u \stackrel{\text{def}}{=} \begin{pmatrix} -A^{-1}v \\ 1 \end{pmatrix}$ , the identity

$$\begin{pmatrix} A & v \\ v^\top & d \end{pmatrix}^\dagger - \begin{pmatrix} A & 0 \\ 0^\top & 0 \end{pmatrix}^\dagger = \frac{1}{d - v^\top A^{-1}v} uu^\top \quad (5)$$

holds.

Note that since the difference in (5) is a rank-one matrix and, therefore, it is positive semi-definite. A recursive application of this property also yields that if  $T \supseteq S$ , then  $W_T^* \succeq W_S^*$ . The second lemma we introduce concerns the inverses of Stieltjes matrices in particular.

**Lemma 2** (Supermodular inverses, Atamtürk and Gómez [3]). *Given a matrix  $Q \in \mathbb{R}^{n \times n}$  and a pair of indices  $i, j \in [n]$ , define the set function  $\theta_{ij}(S) \stackrel{\text{def}}{=} (Q_S^*)_{ij}$  as the  $(i, j)$ -th entry of matrix  $Q_S^*$ . If  $Q$  is a Stieltjes matrix, then  $\theta_{ij}$  is a non-decreasing supermodular function for all pairs of indices  $i, j \in [n]$ .*

Observe that since  $\theta_{ij}(\emptyset) = 0$  and  $\theta_{ij}$  is non-decreasing, Lemma 2 implies that inverses of Stieltjes matrices are non-negative, a fact that is observed in [23] (along with several other properties of Stieltjes matrices).

## 4 Facial structure of Stieltjes polytopes

In this section, we study convexifications of Stieltjes polytopes as defined in Definition 1. Throughout, matrix  $Q \in \mathbb{R}^{n \times n}$  is assumed to be Stieltjes.

**Proposition 1.** Any point  $(\mathbf{z}, \mathbf{W}) \in Z_Q$  satisfies the following properties:

1.  $W_{ij} = W_{ji}$  for all  $1 \leq i < j \leq n$ ,
2.  $\sum_{j=1}^n Q_{ij} W_{ij} = z_i$  for all  $i = 1, \dots, n$ ,
3.  $W_{ij} = 0$  for all  $i, j$  such that  $Q_{ij}^{-1} = 0$ ,
4.  $W_{ij} \geq 0$  for all  $i, j$ .

*Proof.* We first check the validity of the equalities. The first set of equalities follows since all extreme points of  $Z_Q$  have symmetric matrices. The second set follows from [27, Proposition 6]. For the third set, from Lemma 2 (non-negative and non-decreasing  $\theta$ ) it follows that extreme points satisfy  $0 \leq W_{ij} \leq Q_{ij}^{-1}$ : clearly, if  $Q_{ij}^{-1} = 0$ , then  $W_{ij} = 0$  holds. Finally, the fourth set of inequalities follows directly from the non-negativity of inverses of Stieltjes matrices.  $\square$

It is evident from Proposition 1 that  $Z_Q$  is not full-dimensional. In this paper, we study the relaxation induced by the upper-bound constraints

$$P_Q^{\leq} \stackrel{\text{def}}{=} \left\{ (\mathbf{z}, \mathbf{W}) \in \{0, 1\}^n \times \mathbb{R}^{n \times n} : \mathbf{W} \leq (\mathbf{Q} \circ \mathbf{z}\mathbf{z}^\top)^\dagger \right\}.$$

Defining  $Z_Q^{\leq} \stackrel{\text{def}}{=} \text{conv}(P_Q^{\leq})$ , it is easy to show that  $Z_Q^{\leq}$  is full-dimensional.

Given  $i, j \in [n]$  and  $S \subseteq [n]$ , let  $\theta_{ij}(S) = (\mathbf{Q}_S^*)_{ij}$  be the function introduced in Lemma 2, and given  $k \in [n] \setminus S$ , let

$$\rho_{ij}(k; S) \stackrel{\text{def}}{=} \theta_{ij}(S \cup \{k\}) - \theta_{ij}(S).$$

Finally, let  $\mathbf{R}(k; S)$  be the matrix that collects functions  $\rho$  in its entries; that is,  $R(k; S)_{ij} = \rho_{ij}(k; S)$ . Recall, from Corollary 1 and the discussion immediately thereafter, that  $\mathbf{R}$  is a rank-one matrix for all values of  $k$  and  $S$ . Supermodularity of  $\theta_{ij}$  immediately leads to a class of valid inequalities.

#### 4.1 Valid inequalities for $P_Q^{\leq}$

We now study the facial structure of  $Z_Q^{\leq}$ , which bounds matrix  $\mathbf{W}$  from above. The facets of  $Z_Q^{\leq}$  are given by the *polymatroid inequalities* [6], and are a direct consequence of supermodularity of  $\theta$ . For any permutation  $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_n)$  of  $[n]$ , define  $S_0^\pi \stackrel{\text{def}}{=} \emptyset$  and for  $k \in [n]$ , define set

$$S_k^\pi \stackrel{\text{def}}{=} \{\pi_1, \pi_2, \dots, \pi_k\}.$$

**Proposition 2.** For any  $i, j \in [n]$  and any permutation  $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_n)$  of  $[n]$ , the inequality

$$W_{ij} \leq \sum_{k=1}^n \rho_{ij}(\pi_k; S_{k-1}^\pi) z_{\pi_k} \tag{6}$$

is valid for  $Z_Q^{\leq}$ .

*Proof.* Inequality (6) is a polymatroid inequality, necessary to describe the Lovász extension, which describes the concave envelope of a supermodular function [7, 21].  $\square$

Observe that, given permutation  $\pi$ , inequalities (6) can be written compactly (for all  $i, j \in [n]$ ) as the matrix inequalities

$$\mathbf{W} \leq \sum_{k=1}^n \mathbf{R}(\pi_k; S_{k-1}^\pi) z_{\pi_k}. \quad (7)$$

*Example 1 (Continued).* For matrix  $\mathbf{Q} = \begin{pmatrix} 2 & -1 & -1 \\ -1 & 3 & -1 \\ -1 & -1 & 2 \end{pmatrix}$  and permutation  $\pi =$

(1, 2, 3), inequalities (7) reduce to

$$\begin{pmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \end{pmatrix} \leq \begin{pmatrix} 1/2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} z_1 + \begin{pmatrix} 1/10 & 1/5 & 0 \\ 1/5 & 2/5 & 0 \\ 0 & 0 & 0 \end{pmatrix} z_2 + \begin{pmatrix} 16/15 & 4/5 & 4/3 \\ 4/5 & 3/5 & 1 \\ 4/3 & 1 & 5/3 \end{pmatrix} z_3. \blacksquare$$

## 4.2 Strength of the inequalities

We now show that inequalities (7) are indeed strong.

**Proposition 3.** *Inequality (6) is facet-defining for  $\text{conv}(P_Q^\leq)$ .*

*Proof.* In Table 1, we provide  $(n + n^2)$  affinely independent points in  $P_Q^\leq$  such that (6) holds at equality.

Table 1: Affinely independent points in  $P_Q^\leq$  satisfying (6) at equality.

| # | Indices   | $z$                    | $\mathbf{W}$                           |
|---|---|------------------------|--|
| 1 | -   | $\mathbf{e}$           | $\mathbf{Q}^{-1}$                      |
| 2 | $\forall k, \ell \in [n]$ with $(k, \ell) \neq (i, j)$ and $Q_{ij}^{-1} \neq 0$ | $\mathbf{e}$           | $\mathbf{Q}^{-1} - \mathbf{E}_{k\ell}$ |
| 3 | $\forall k \in [n]$   | $\mathbf{e}_{S_k^\pi}$ | $\mathbf{Q}_{S_k^\pi}^*$               |

- Point #1 belongs to  $P_Q$  and, therefore, it also belongs to  $P_Q^\leq$ ; if  $z = \mathbf{e}$ , then inequality (6) reduces to  $W_{ij} \leq \theta_{ij}(\mathbf{e}) = Q_{ij}^{-1}$ , and thus the inequality is satisfied at equality at this point.

- Points #2 correspond to  $n^2 - 1$  points which are obtained by subtracting a non-negative quantity from the feasible point #1 and, therefore, also belong to  $P_Q^\leq$ ; the inequality is tight for the same reason as point #1. Points #2 are affinely independent from previously introduced points because they are the first points with  $W_{k\ell} \neq Q_{k\ell}^{-1}$ .

- Points #3 correspond to  $n$  points belonging to  $P_Q$ ; therefore, they also belong to  $P_Q^\leq$ ; if  $z = \mathbf{e}_{S_k^\pi}$ , then inequality (6) reduces to  $W_{ij} \leq \theta_{ij}(\mathbf{e}_{S_k^\pi}) = \left( \mathbf{Q}_{S_k^\pi}^* \right)_{ij}$ , and thus the inequality is satisfied at equality at this point. Finally,

points #3 are affinely independent from others because they are the first points with  $z_k \neq 1$ .  $\square$

From Proposition 3, we see that all inequalities (6) for all combinations of  $i, j \in [n]$  and all permutations  $\pi$  are necessary to describe  $Z_Q^\leq$ . We now show that they are, along with the bound inequalities, sufficient.

**Theorem 2.** *Inequalities (7) (for all permutations  $\pi$ ) and bound constraints  $\mathbf{0} \leq \mathbf{z} \leq \mathbf{e}$  describe  $\text{conv}(P_Q^\leq)$ .*

*Proof.* We show that, for any  $\mathbf{c} \in \mathbb{R}^n$  and  $\Sigma \in \mathbb{R}^{n \times n}$ , the optimization problems

$$\min_{\mathbf{z}, \mathbf{W}} \mathbf{c}^\top \mathbf{z} + \langle \Sigma, \mathbf{W} \rangle \text{ s.t. } (\mathbf{z}, \mathbf{W}) \in P_Q^\leq \quad (8)$$

$$\min_{\mathbf{z}, \mathbf{W}} \mathbf{c}^\top \mathbf{z} + \langle \Sigma, \mathbf{W} \rangle \text{ s.t. } (7), \mathbf{0} \leq \mathbf{z} \leq \mathbf{e} \quad (9)$$

are equivalent; that is, either both are unbounded, or there exists an optimal solution of (9) that is feasible for (8).

Note that if  $\Sigma_{ij} > 0$  for any  $i, j \in [n]$ , then both problems are unbounded by letting  $W_{ij} \rightarrow -\infty$ . Therefore, we assume  $\Sigma \leq \mathbf{0}$ . In this case, in optimal solutions of (8),  $\mathbf{W}$  will be set to its upper bound, i.e.,  $\mathbf{W} = (\mathbf{Q} \circ \mathbf{z}\mathbf{z}^\top)^\dagger$  holds. By abuse of notation, let  $\theta_{ij}(\mathbf{z}) \stackrel{\text{def}}{=} \theta_{ij}(S_z)$ , where  $S_z = \{i \in [n] : z_i = 1\}$ ; since  $((\mathbf{Q} \circ \mathbf{z}\mathbf{z}^\top)^\dagger)_{ij} = \theta_{ij}(\mathbf{z})$ , we find that (8) is equivalent to

$$\min_{\mathbf{z} \in \{0,1\}^n} \mathbf{c}^\top \mathbf{z} + \sum_{i=1}^n \sum_{j=1}^n \Sigma_{ij} \theta_{ij}(\mathbf{z}). \quad (10)$$

Recalling that functions  $\theta_{ij}$  are supermodular (Lemma 2) and  $\Sigma_{ij} \leq 0$  for all  $i, j$ , we find that  $\Theta(\mathbf{z}) \stackrel{\text{def}}{=} \sum_{i=1}^n \sum_{j=1}^n \Sigma_{ij} \theta_{ij}(\mathbf{z})$  is a submodular function. Therefore, it follows that (10) is equivalent to minimization over its Lovász extension [21], which is precisely (9).  $\square$

Now consider the relaxation of (1) induced by Theorem 1, but using only inequalities (7) and bound constraints instead of the full description of  $P_Q$ :

$$\min_{\mathbf{x}, \mathbf{z}, \mathbf{W}, t} \mathbf{a}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{z} + t \quad (11a)$$

$$\text{s.t. } \begin{pmatrix} \mathbf{W} & \mathbf{x} \\ \mathbf{x}^\top & t \end{pmatrix} \in \mathbb{S}_+^{n+1} \quad (11b)$$

$$\mathbf{W} \leq \sum_{k=1}^n \mathbf{R}(\pi_k; S_{k-1}^\pi) z_{\pi_k} \quad \text{for all permutations } \pi \text{ of } [n] \quad (11c)$$

$$(\mathbf{x}, \mathbf{z}) \in C \quad (11d)$$

$$\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in [0, 1]^n, t \in \mathbb{R}_+, \mathbf{W} \in \mathbb{R}^{n \times n}. \quad (11e)$$

Because constraints (11c) are a relaxation of constraints  $(z, \mathbf{W}) \in P_Q$ , we find from Theorem 1 that (11) is indeed a valid relaxation. In general, (11) can be weak: in fact, it can be unbounded unless constraints  $\mathbf{W} \geq \mathbf{0}$  are also added. Nonetheless, as we now show, under the specific conditions stated in [3] for polynomial-time solvability of (1), the relaxation (11) is exact. Given a Stieltjes matrix  $\mathbf{Q}$ , define the optimization

$$\min_{\substack{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n \\ \mathbf{x} \circ (\mathbf{e} - \mathbf{z}) = \mathbf{0}}} \mathbf{a}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{z} + \mathbf{x}^\top \mathbf{Q} \mathbf{x}, \quad (12)$$

which is the special case of (1) with  $C = \mathbb{R}^{2n}$ .

**Proposition 4.** *If  $\mathbf{a} \leq \mathbf{0}$  or  $\mathbf{a} \geq \mathbf{0}$ , and  $C = \mathbb{R}^{2n}$ , then there exists an optimal solution of (11) that is also optimal for (12), with the same objective value.*

*Proof.* The start of the proof follows the steps in [27, Theorem 1], which we repeat for completeness. Constraint (11b) is equivalent to the system [2]

$$\mathbf{W} \succeq 0, \quad t \geq \mathbf{x}^\top \mathbf{W}^\dagger \mathbf{x}, \quad \text{and} \quad \mathbf{W} \mathbf{W}^\dagger \mathbf{x} = \mathbf{x}.$$

Therefore, variable  $t$  can be easily projected out since any optimal solution satisfies  $t = \mathbf{x}^\top \mathbf{W}^\dagger \mathbf{x}$ . We can restate problem (11) as

$$\min_{\mathbf{x}, \mathbf{z}, \mathbf{W}} \mathbf{a}^\top \mathbf{W} \mathbf{W}^\dagger \mathbf{x} + \mathbf{c}^\top \mathbf{z} + \mathbf{x}^\top \mathbf{W}^\dagger \mathbf{x} \quad (13a)$$

$$\text{s.t. } \mathbf{W} \mathbf{W}^\dagger \mathbf{x} = \mathbf{x}, \quad \mathbf{W} \in \mathbb{S}_+^n \quad (13b)$$

$$(11c) - (11e). \quad (13c)$$

Note that we use the equality in (13b) to rewrite a linear term in the objective. We now project out variables  $\mathbf{x}$ : The KKT conditions associated with the continuous variables are

$$\begin{aligned} \mathbf{W} \mathbf{W}^\dagger \mathbf{x} &= \mathbf{x} \\ \mathbf{a}^\top \mathbf{W} \mathbf{W}^\dagger + 2\mathbf{W}^\dagger \mathbf{x} + \boldsymbol{\lambda}^\top (\mathbf{W} \mathbf{W}^\dagger - \mathbf{I}) &= \mathbf{0}, \end{aligned}$$

which are satisfied by setting  $\mathbf{x}^* = -\frac{1}{2}\mathbf{W}\mathbf{a}$  and  $\boldsymbol{\lambda}^* = \mathbf{0}$ . Substituting  $\mathbf{x}$  with its optimal value, we find that problem (13) further simplifies to

$$\min_{\mathbf{z}, \mathbf{W}} \left\langle -\frac{1}{4}\mathbf{a}\mathbf{a}^\top, \mathbf{W} \right\rangle + \mathbf{c}^\top \mathbf{z} \quad (14a)$$

$$\text{s.t. } (11c), \quad \mathbf{W} \in \mathbb{S}_+^n, \quad \mathbf{0} \leq \mathbf{z} \leq \mathbf{e}. \quad (14b)$$

In the last step of the proof (which does not follow from [27]), we study the relaxation of (14) obtained by removing constraint  $\mathbf{W} \in \mathbb{S}_+^n$ . In particular, we show that there exist optimal solutions  $(\bar{\mathbf{z}}, \bar{\mathbf{W}}, \bar{t})$  of the relaxation such that: (i)  $\bar{\mathbf{z}}$  is integral; (ii)  $\bar{\mathbf{W}}$  is positive semidefinite; (iii)  $\bar{t} = -\frac{1}{4}\mathbf{a}_S^\top \mathbf{Q}_S^{-1} \mathbf{a}_S$ , where  $S = \{i \in [n] : \bar{z}_i = 1\}$ . As a consequence,  $(\bar{\mathbf{z}}, \bar{\mathbf{W}}, \bar{t})$  is also optimal for (14) and (12), concluding the proof.

Observe that  $-\frac{1}{4}\mathbf{a}\mathbf{a}^\top \leq \mathbf{0}$  due to the assumption that  $\mathbf{a}$  is of the same sign. Therefore, if  $\mathbf{W} \in \mathbb{S}_+^n$  is removed, then in optimal solutions of the relaxation,  $\bar{\mathbf{W}}$  is equal to its upper bound, i.e.,

$$\bar{\mathbf{W}} = \min_{\pi \in \Pi} \sum_{k=1}^n \mathbf{R}(\pi_k; S_{k-1}^\pi) z_{\pi_k}, \quad (15)$$

where  $\Pi$  is the set of all permutations of  $[n]$ . In other words, the relaxation is equivalent to

$$\min_{\mathbf{0} \leq \mathbf{z} \leq \mathbf{1}} \min_{\pi \in \Pi} \left\{ \sum_{k=1}^n \left( \left\langle -\frac{1}{4}\mathbf{a}\mathbf{a}^\top, \mathbf{R}(\pi_k; S_{k-1}^\pi) \right\rangle z_{\pi_k} \right) \right\} + \mathbf{c}^\top \mathbf{z}. \quad (16)$$

Recognizing (16) as a linear optimization problem over the Lovász extension of a submodular function, we conclude that  $\bar{\mathbf{z}} \in \{0, 1\}^n$  in optimal solutions, proving (i). Since optimal permutations for (15) correspond to nondecreasing orders of  $\bar{\mathbf{z}}$  [10], letting  $\tau = \|\bar{\mathbf{z}}\|_0$  we conclude that

$$\bar{\mathbf{W}} = \sum_{k=1}^{\tau} \mathbf{R}(\pi_k; S_{k-1}) = \sum_{k=1}^{\tau} \left( \mathbf{W}_{S_k}^* - \mathbf{W}_{S_{k-1}}^* \right) = \mathbf{W}_{S_\tau}^*,$$

where we defined  $\mathbf{W}_{S_0}^* = \mathbf{0}$ . In particular,  $\bar{\mathbf{W}} \in \mathbb{S}_+^n$ , proving (ii), and substituting  $\mathbf{W}$  with its optimal value in (16) we can prove (iii).  $\square$

### 4.3 Separation algorithm

We now discuss the separation problem: given any fixed point  $(\bar{\mathbf{z}}, \bar{\mathbf{W}}) \in \mathbb{R}^n \times \mathbb{R}^{n \times n}$ , how to find a violated inequality (6) if there exists one. Since inequality (6) is a polymatroid inequality, it follows that the most violated inequalities correspond to permutations  $(\pi_1, \pi_2, \dots, \pi_n)$  such that  $\bar{z}_{\pi_1} \geq \bar{z}_{\pi_2} \geq \dots \geq \bar{z}_{\pi_n}$  [10]. In particular, since the permutation depends only on the values of  $\bar{\mathbf{z}}$ , we find that most violated permutations coincide for all values of indices  $i, j$  in (6). Therefore, using the matrix notation (7), we see that a straightforward way to compute a violated inequality is simply to compute matrices  $\mathbf{Q}_{S_k}^*$  for all  $k \in [n]$ , which requires inverting  $\mathcal{O}(n)$  matrices. A faster approach to compute the coefficients of inequalities (7) consists of using a Cholesky decomposition.

Indeed, consider the properties of coefficient matrices  $\{\mathbf{R}(\pi_k; S_{k-1})\}_{k=1}^n$  in inequalities (7). First, they are rank-one matrices; that is, there exists  $\mathbf{v}_k \in \mathbb{R}^n$  such that  $\mathbf{R}(\pi_k; S_{k-1}) = \mathbf{v}_k \mathbf{v}_k^\top$ . Second, they add up to  $\mathbf{Q}^{-1}$ ; that is,  $\sum_{k=1}^n \mathbf{R}(\pi_k; S_{k-1}) = \sum_{k=1}^n \mathbf{v}_k \mathbf{v}_k^\top = \mathbf{Q}^{-1}$ . Third, the entries corresponding to indices that have not appeared yet in the permutation vanish; that is,  $R(\pi_k; S_{k-1})_{ij} = 0$  if  $\max\{i, j\} > k$ , or equivalently  $(v_k)_i = 0$  if  $i > k$ . If we define a matrix  $\mathbf{V} \in \mathbb{R}^{n \times n}$  such that its  $i$ -th column is precisely  $\mathbf{v}_i$ , these properties are equivalent to  $\mathbf{V}\mathbf{V}^\top = \mathbf{Q}^{-1}$  and  $\mathbf{V}$  is upper triangular (where elements are ordered according to the permutation  $\pi$ ). Note that these properties

are similar to the ones corresponding to a Cholesky decomposition of  $\mathbf{Q}^{-1}$ , except that the matrix in the Cholesky decomposition is lower triangular instead of upper triangular. To account for this difference, it suffices to compute the Cholesky decomposition in the *reverse order* of the permutation, and because  $\mathbf{Q}^{-1} \in \mathbb{S}_{++}^n$ , the Cholesky decomposition is the unique matrix satisfying the required properties, and thus indeed coincides with the coefficients in inequalities (7). We summarize the separation algorithm in Proposition 5 below.

**Proposition 5.** *Algorithm 1 produces up to  $\mathcal{O}(n^2)$  violated inequalities, one for each combination of elements  $i, j \in [n]$ , if there exists any.*

---

**Algorithm 1** Separation procedure

---

**Input:** Point  $\bar{\mathbf{z}} \in \mathbb{R}^n$ .

**Output:** Most violated inequalities.

- 1: Find a permutation  $\pi$  satisfying  $\bar{z}_{\pi_1} \leq \bar{z}_{\pi_2} \leq \dots \leq \bar{z}_{\pi_n}$  ▷ Sorting
- 2: Compute Cholesky decomposition  $\mathbf{Q}^{-1} = \mathbf{V}\mathbf{V}^\top$  according to order  $\pi$
- 3: **return** inequalities

$$\mathbf{W} \leq \sum_{k=1}^n (\mathbf{v}_{\pi_k} \mathbf{v}_{\pi_k}^\top) z_{\pi_k}$$

where  $\mathbf{v}_{\pi_i}$  denotes the  $i$ -th column of  $\mathbf{V}$ .

---

We emphasize that the order in line 1 of Algorithm 1 is non-decreasing, instead of the natural non-increasing order that arises often with polymatroid inequalities, since we use the *reverse order* of the permutation in computing the Cholesky decomposition as we discussed in the preceding paragraph.

We now discuss the runtime of Algorithm 1. Since sorting the variables (line 1) can be done in  $\mathcal{O}(n \log n)$ , we find that the complexity of Algorithm 1 is dominated by the cost of computing a Cholesky decomposition (line 2), which is usually  $\mathcal{O}(n^3)$ . Note that line 2 also requires inverting matrix  $\mathbf{Q}$ , but this operation (with cubic runtime as well) needs to be performed *only* once as preprocessing. In some cases, the runtime for a general Stieltjes matrix can be improved. For example, Lemma 1 can be called recursively to compute the coefficients, requiring  $\mathcal{O}(n)$  calls to a routine for matrix-vector product and matrix-matrix subtraction: if  $\mathbf{Q}$  is sparse, then the vectors appearing in the computations are sparse as well, potentially improving runtimes for the multiplications (with appropriate data structures). Finally, we point out that the output of Algorithm 1 consists of  $\mathcal{O}(n^2)$  inequalities with up to  $n$  nonzeros per inequalities, thus formulating the inequalities in a solver already requires processing a cubic number of nonzeros, thus it is not possible to improve this runtime (at least with an off-the-shelf solver).

## 5 Computational results

In this section, we describe computational experiments performed to test the effectiveness of the proposed convexification. For the computational study, we consider problems of the form (3), arising from the inference of sparse Besag-York-Mollie graphical models as discussed in §2.1. Specifically, given a graph  $\mathcal{G} = (V, E)$ , we consider problem

$$\min_{\mathbf{x} \in \mathbb{R}^V, \mathbf{z} \in \{0,1\}^V} \frac{1}{\sigma^2} \sum_{i \in V} (y_i - x_i)^2 + \sum_{[i,j] \in E} (x_i - x_j)^2 + \mu \sum_{i \in V} z_i \quad (17a)$$

$$\text{s.t.} \quad \sum_{i \in V} z_i \leq k, \quad \mathbf{x} \circ (\mathbf{e} - \mathbf{z}) = \mathbf{0}, \quad (17b)$$

where  $\sigma, \mu, k \in \mathbb{R}_+$  are given parameters. In our experiments, we solve *cardinality-penalized problems* with  $\mu > 0$  and  $k = |V|$  (modeling situations where the density of the model is penalized) as well as *cardinality-constrained problems* with  $\mu = 0$  and  $k < |V|$  (modeling situations where the sparsity is directly specified). In either case, we solve (17) for different values of  $\sigma$ , a parameter that is connected to the noise of the underlying model.

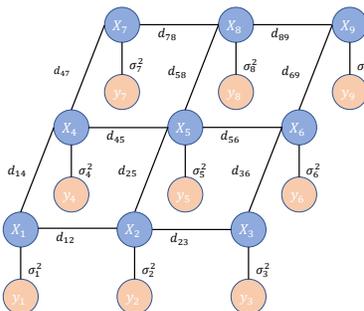


Figure 1: Grid graph for modeling the Besag-York-Mollie process; from [19].

The data are generated following [19] and are available online at <https://sites.google.com/usc.edu/gomez/data>. Figure 1 depicts the graph  $\mathcal{G}$  used to model spatial processes. Graph  $\mathcal{G}$  is given by a two-dimensional lattice; that is, vertices  $V = [n]$  are arranged in a grid, with edges between horizontally and vertically adjacent vertices. We consider instances with grid sizes  $10 \times 10$ , thus resulting in instances with  $|V| = 100$ . To generate the data, we create a sparse “true” signal  $\{X_i\}_{i \in V}$  with three non-negative spikes, each spike affecting a  $3 \times 3$  block of elements in  $V$ . Thus, the generated signal can have at most 27 non-zeros (the construction of the underlying true signal is explained in detail in Appendix A). We then generate noisy observations as  $y_i = |X_i + \epsilon_i|$ , where  $\epsilon_i$  are i.i.d. samples from a Gaussian distribution  $\mathcal{N}(0, \sigma^2)$ . Note that since  $\mathbf{y} \geq \mathbf{0}$ ,

linear coefficients of variables  $\mathbf{x}$  obtained by expanding the quadratic terms in (17a) are non-positive.

In § 5.1, we describe the models tested, and in § 5.2, we report the computational experiments.

## 5.1 Models

We compare three reformulations or relaxations of the problem (1). The former two are commonly used perspective formulations [11, 14], while the latter is based on the valid inequalities proposed in §4. We describe these formulations next.

**pers-c:** Perspective relaxation with continuous variable  $\mathbf{z}$  and big-M constraints:

$$\min_{\mathbf{x}, \mathbf{z}} \frac{1}{\sigma^2} \|\mathbf{y}\|_2^2 + \frac{1}{\sigma^2} \sum_{i \in V} \left( \frac{x_i^2}{z_i} - 2y_i x_i \right) + \sum_{[i,j] \in E} (x_i - x_j)^2 + \mu \sum_{i \in V} z_i \quad (18a)$$

$$\text{s.t.} \quad \sum_{i \in V} z_i \leq k, \quad -M\mathbf{z} \leq \mathbf{x} \leq M\mathbf{z} \quad (18b)$$

$$\mathbf{x} \in \mathbb{R}^V, \mathbf{z} \in [0, 1]^V, \quad (18c)$$

where the value  $M$  is large enough so that the big-M constraints are redundant when  $z_i = 1$ . In our computations, the bounds are set to  $M = 10$ .

**pers-b:** Perspective reformulation **pers-c** with binary variables  $\mathbf{z} \in \{0, 1\}^V$ . Using this formulation, problems are solved to optimality with branch-and-bound, closing the gap from the perspective relaxation.

**poly:** Polymatroid relaxation with the equalities and lower bounds given in Proposition 1 and the polymatroid inequalities as described in (11):

$$\min_{\mathbf{x}, \mathbf{z}, \mathbf{W}, t} \frac{1}{\sigma^2} \|\mathbf{y}\|_2^2 - 2 \frac{1}{\sigma^2} \sum_{i \in V} y_i x_i + \mu \sum_{i \in V} z_i + t \quad (19a)$$

$$\text{s.t.} \quad \sum_{i \in V} z_i \leq k \quad (19b)$$

$$\begin{pmatrix} \mathbf{W} & \mathbf{x} \\ \mathbf{x}^\top & t \end{pmatrix} \succeq 0 \quad (19c)$$

$$W_{ij} = W_{ji} \quad \text{for all } i, j \in V \quad (19d)$$

$$\sum_{j \in V} Q_{ij} W_{ij} = z_i \quad \text{for all } i \in V \quad (19e)$$

$$\mathbf{W} \leq \sum_{k \in V} \mathbf{R}(\pi_k; S_{k-1}^\pi) z_{\pi_k} \quad \text{for all permutations } \pi \text{ of } V \quad (19f)$$

$$\mathbf{x} \in \mathbb{R}^V, \mathbf{z} \in [0, 1]^V, t \in \mathbb{R}_+, \mathbf{W} \in \mathbb{R}_+^{V \times V}, \quad (19g)$$

where  $\mathbf{Q}$  is the Stieltjes matrix obtained by collecting all nonlinear terms in the objective and inequalities (19f) are those required to describe  $\text{conv}(P_{\mathbf{Q}}^{\leq})$ . Note that (19) is a reformulation of (1) under the conditions of Proposition 4, that is, in the cardinality-penalized instances, and is a relaxation otherwise.

*Remark 1* (Implementation of `poly`). As the total number of polymatroid inequalities is  $|V|!$ , listing all of them in (19) is impractical. Instead, we add (19f) as cutting planes; that is, we solve the `poly` problem (19) by adding polymatroid inequalities (19f) as cutting planes iteratively using the separation procedure 1 and stopping when the difference of the objective values between two subsequent iterations is small enough ( $< 10^{-3}$ ).

## 5.2 Results

The perspective formulations (`pers-c`, `pers-b`) are solved with Gurobi version 9.0.2 using a single thread. The time limit for the computations is set to one hour; all other configurations are set to default values. The polymatroid relaxation (`poly`) is solved with Mosek version 9.3 with the default settings. All experiments are run on a Lenovo laptop with a 1.9 GHz Intel®Core™ i7-8650U CPU and 16 GB main memory. For cardinality-constrained instances, we set  $k = 20$ , and for cardinality-penalized instances, we choose  $\mu$  so that the number of non-zeros of the estimator approximately matches the number of non-zeros of the underlying signal.

Tables 2 and 3 present results for values of  $\sigma^2 \in \{0.5, 1.0, 2.0, 5.0, 10.0\}$  and cardinality-penalized and cardinality-constrained instances, respectively. Each row represents the average over five instances generated with identical parameters. The gap corresponding to each model compares the lower bound produced by the model with the best upper bound found by Gurobi with model `pers-b`. For model `pers-b`, we also report, under the `#Opt` column, the number of instances that can be solved to optimality before the time limit, and, under the `Gap` column, we report the average optimality gap at termination. The number of iterations for model `poly` is the number of cutting plane rounds required before termination, with each iteration resulting in the addition of up to  $\binom{|V|}{2}$  cuts. All times reported are in seconds.

Table 2: Experiments with cardinality-penalized instances.

| $\sigma^2$ | $\mu$ | pers-c |      | pers-b |         |      | poly  |       |                   |
|------------|-------|--------|------|--------|---------|------|-------|-------|-------------------|
|            |       | Time   | Gap  | #Opt   | Time    | Gap  | #Iter | Time  | Gap               |
| 0.5        | 0.25  | 0.1    | 2.4% | 5      | 2.1     | 0.0% | 4     | 102.4 | $3 \cdot 10^{-8}$ |
| 1.0        | 0.12  | 0.1    | 5.7% | 1      | 2,998.3 | 1.7% | 6     | 235.4 | $7 \cdot 10^{-4}$ |
| 2.0        | 0.12  | 0.1    | 5.1% | 2      | 2,937.1 | 0.5% | 6     | 218.3 | $2 \cdot 10^{-7}$ |
| 5.0        | 0.12  | 0.1    | 5.1% | 1      | 3,021.9 | 0.9% | 6     | 235.6 | $3 \cdot 10^{-6}$ |
| 10.0       | 0.12  | 0.1    | 5.1% | 2      | 2,990.6 | 0.8% | 6     | 219.0 | $9 \cdot 10^{-8}$ |

Table 3: Experiments with cardinality-constrained instances.

| $\sigma^2$ | poly |      | pers-b |         |      | poly  |       |                   |
|------------|------|------|--------|---------|------|-------|-------|-------------------|
|            | Time | Gap  | #Opt   | Time    | Gap  | #Iter | Time  | Gap               |
| 0.5        | 0.1  | 3.3% | 5      | 13.1    | 0.0% | 4     | 181.3 | $4 \cdot 10^{-7}$ |
| 1.0        | 0.1  | 7.4% | 0      | 3,600.0 | 1.5% | 7     | 285.3 | $9 \cdot 10^{-4}$ |
| 2.0        | 0.1  | 7.1% | 1      | 3,017.3 | 1.3% | 7     | 277.5 | $3 \cdot 10^{-3}$ |
| 5.0        | 0.1  | 7.1% | 1      | 2,990.3 | 1.4% | 7     | 281.7 | $2 \cdot 10^{-3}$ |
| 10.0       | 0.1  | 7.1% | 1      | 2,988.9 | 1.3% | 7     | 260.0 | $3 \cdot 10^{-3}$ |

In both cases, we observe that the relaxations from the perspective relaxation can be solved in a fraction of a second and yield gaps between 2% and 7%: the gaps are larger for cardinality-constrained instances and for problems with larger values of  $\sigma$ . Indeed, for a larger value of  $\sigma$ , terms  $x_i^2/z_i$  (crucial to the strength of the perspective relaxation) in the objective of (18) represent a relatively smaller portion of the objective, leading to larger gaps. When used in a branch-and-bound algorithm, it is possible to efficiently prove optimality if  $\sigma^2 = 0.5$ , but most instances with larger values of  $\sigma$  cannot be solved within the time limit of one hour.

For cardinality-penalized instances, model **poly** delivers optimal solutions on average in about 200 seconds. The runtimes are larger than those required to solve the perspective relaxation **pers-c** due to expenses associated with solving a semidefinite program with a large number of linear inequalities via cutting planes. Nonetheless, for the more challenging instances with  $\sigma^2 \geq 1$ , the runtimes of **poly** are at least an order of magnitude less than those arising from the branch-and-bound method to solve **pers-b**; for these instances, both **poly** and **pers-b** are exact models. The runtimes of **poly** for cardinality-constrained instances are similar. Although in this case, due to the additional cardinality constraint, **poly** is not guaranteed to deliver exact solutions to problem (17), the gaps proved are almost 0% and much smaller than those obtained after one hour of branching with **pers-b**.

## Acknowledgments

Alper Atamtürk is supported, in part, by NSF AI Institute for Advances in Optimization Award 2112533 and the Office of Naval Research grant N00014-24-1-2149. Andrés Gómez is supported, in part, by grant FA9550-22-1-0369 from the Air Force Office of Scientific Research. Simge Küçükyavuz is supported, in part, by grant #2007814 from the National Science Foundation.

## References

- [1] M. S. Aktürk, A. Atamtürk, and S. Gürel. A strong conic quadratic reformulation for machine-job assignment with controllable processing times. *Operations Research Letters*, 37:187–191, 2009.
- [2] A. Albert. Conditions for positive and nonnegative definiteness in terms of pseudoinverses. *SIAM Journal on Applied Mathematics*, 17(2):434–440, 1969.
- [3] A. Atamtürk and A. Gómez. Strong formulations for quadratic optimization with M-matrices and indicator variables. *Mathematical Programming*, 170:141–176, 2018.
- [4] A. Atamtürk and A. Gómez. Rank-one convexification for sparse regression. *arXiv preprint arXiv:1901.10334*, 2019.
- [5] A. Atamtürk and A. Gómez. Supermodularity and valid inequalities for quadratic optimization with indicators. *Mathematical Programming*, 201:295–338, 2023.
- [6] A. Atamtürk and V. Narayanan. Polymatroids and mean-risk minimization in discrete optimization. *Operations Research Letters*, 36:618–622, 2008.
- [7] A. Atamtürk and V. Narayanan. Submodular function minimization and polarity. *Mathematical Programming*, 196:57–67, 2022.
- [8] A. Atamtürk, A. Gómez, and S. Han. Sparse and smooth signal estimation: Convexification of  $\ell_0$  formulations. *Journal of Machine Learning Research*, 22:1–43, 2021.
- [9] J. Besag, J. York, and A. Mollié. Bayesian image restoration, with two applications in spatial statistics. *Annals of the Institute of Statistical Mathematics*, 43(1):1–20, 1991.
- [10] J. Edmonds. Submodular functions, matroids, and certain polyhedra. In R. Guy, H. Hanani, N. Sauer, and J. Schönheim, editors, *Combinatorial Structures and Their Applications*, pages 69–87. Gordon and Breach, 1970.
- [11] A. Frangioni and C. Gentile. Perspective cuts for a class of convex 0–1 mixed integer programs. *Mathematical Programming*, 106:225–236, 2006.
- [12] A. Frangioni and C. Gentile. SDP diagonalizations and perspective cuts for a class of nonseparable MIQP. *Operations Research Letters*, 35:181–185, 2007.
- [13] A. Frangioni, C. Gentile, and J. Hungerford. Decompositions of semidefinite matrices and the perspective reformulation of nonseparable quadratic programs. *Mathematics of Operations Research*, 45(1):15–33, 2020.

- [14] O. Günlük and J. Linderoth. Perspective reformulations of mixed integer nonlinear programs with indicator variables. *Mathematical Programming*, 124:183–205, 2010.
- [15] S. Han and A. Gómez. Compact extended formulations for low-rank functions with indicator variables. *arXiv preprint arXiv:2110.14884*, 2021.
- [16] S. Han, A. Gómez, and A. Atamtürk.  $2 \times 2$  convexifications for convex quadratic optimization with indicator variables. *arXiv preprint arXiv:2004.07448*, 2020.
- [17] S. Han, A. Gómez, and J.-S. Pang. On polynomial-time solvability of combinatorial Markov random fields. *arXiv preprint arXiv:2209.13161*, 2022.
- [18] S. Han, A. Gómez, and A. Atamtürk.  $2 \times 2$ -convexifications for convex quadratic optimization with indicator variables. *Mathematical Programming*, 202:95–134, 2023.
- [19] Z. He, S. Han, A. Gómez, Y. Cui, and J.-S. Pang. Comparing solution paths of sparse quadratic minimization with a Stieltjes matrix. *Mathematical Programming*, 204:517–566, 2024.
- [20] P. Liu, S. Fattahi, A. Gómez, and S. Küçükyavuz. A graph-based decomposition method for convex quadratic optimization with indicators. *Mathematical Programming*, 200(2):669–701, 2023.
- [21] L. Lovász. Submodular functions and convexity. In *Mathematical programming the state of the art*, pages 235–257. Springer, 1983.
- [22] T.-T. Lu and S.-H. Shiou. Inverses of  $2 \times 2$  block matrices. *Computers & Mathematics with Applications*, 43(1-2):119–129, 2002.
- [23] R. J. Plemmons. M-matrix characterizations. I – nonsingular M-matrices. *Linear Algebra and its Applications*, 18:175–188, 1977.
- [24] S. Shafieezadeh-Abadeh and F. Kılınç-Karzan. Constrained optimization of rank-one functions with indicator variables. *Mathematical Programming*, 2024. doi: <https://doi.org/10.1007/s10107-023-02047-y>. Article in Advance.
- [25] L. Wei, A. Gómez, and S. Küçükyavuz. On the convexification of constrained quadratic optimization problems with indicator variables. In *International Conference on Integer Programming and Combinatorial Optimization*, pages 433–447. Springer, 2020.
- [26] L. Wei, A. Gómez, and S. Küçükyavuz. Ideal formulations for constrained convex optimization problems with indicator variables. *Mathematical Programming*, 192(1-2):57–88, 2022.

- [27] L. Wei, A. Atamtürk, A. Gómez, and S. Küçükyavuz. On the convex hull of convex quadratic optimization problems with indicators. *Mathematical Programming*, 204(1-2):703–737, 2024.

## A Construction of true signals in computational experiments

Here we discuss how to construct the true signal  $\{X_i\}_{i \in V}$ . Note that since each node in  $\mathcal{G}$  is uniquely determined by its coordinates  $1 \leq k, \ell \leq m$  in the  $m \times m$  grid, we let  $X_{(k,\ell)}$  denote the value of the signal at those coordinates.

The true signal is initially set to 0 at all its coordinates. We then repeat the following process three times:

1. Uniformly pick coordinates  $2 \leq k, \ell \leq 9$ .
2. Generate a 9-dimensional Gaussian spike  $\mathbf{s} \sim \mathcal{N}(\mathbf{0}, \Theta^{-1})$ , where

$$\Theta = \begin{pmatrix} 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 4 & 0 & 0 & -1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 4 & -1 & 0 & -1 & 0 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & 0 & -1 & 0 & -1 & 4 & 0 & 0 & -1 \\ 0 & 0 & 0 & -1 & 0 & 0 & 4 & -1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 & -1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 \end{pmatrix}.$$

3. For all  $0 \leq j_1, j_2 \leq 2$ , let  $h = 1 + 3j_1 + j_2$  and update

$$X_{(k-1+j_1, \ell-1+j_2)} \leftarrow X_{(k-1+j_1, \ell-1+j_2)} + |s_h|.$$