

Optimization of Sums of Bivariate Functions

An Introduction to Relaxation-Based Methods
for the Case of Finite Domains

A PREPRINT

Nils Müller

Max Planck Institute for Informatics
Saarbrücken, Germany;
Department of Mathematics
Ruhr University Bochum
Bochum, Germany
nmueller@mpi-inf.mpg.de

November 26, 2025

Abstract. We study the optimization of functions with $n > 2$ arguments that have a representation as a sum of several functions that have only 2 of the n arguments each, termed *sums of bivariates*, on finite domains.

The complexity of optimizing sums of bivariates is shown to be *NP-equivalent* and it is shown that there exists *free lunch* in the optimization of sums of bivariates.

Based on measure-valued extensions of the objective function, so-called *relaxations*, ℓ^2 -approximation, and entropy-regularization, we derive several tractable problem formulations solvable with linear programming, coordinate ascent as well as with closed-form solutions.

The limits of applying tractable versions of such relaxations to sums of bivariates are investigated using general results for reconstructing measures from their bivariate marginals.

Experiments in which the derived algorithms are applied to random functions, vertex coloring, and signal reconstruction problems provide insights into qualitatively different function classes that can be modeled as sums of bivariates.

Keywords Linear Programming · Graphical Models · Relaxation · Inverse Problems

Contents

1. Introduction	3
2. Approximation	7
3. Optimization	12
3.1. Fundamental Perspective	12
3.2. Relaxation Principle	17
3.3. Entropy-Regularized Relaxation Principle	25
3.4. Recovering Primal Solutions from Dual Solutions	30
4. Algorithms	33
4.1. Coordinate Descent for the Sum of Bivariates	33
4.2. Linear Programming for the Dual Linear Program	33
4.3. Block Coordinate Ascent for the Dual Tree Relaxation	36
4.4. Block Coordinate Ascent for the Dual Entropy-Regularized Tree Relaxation	36
5. Experiments	40
5.1. Random Sums of Bivariates	40
5.2. Vertex Coloring	41
5.3. Signal Reconstruction	42
6. Conclusion and Future Work	45
References	46
A. Additional Proofs	48
A.1. Lemmata for Theorem 3.5	48
B. Additional Code	55
B.1. Verifying that $F \circ \tilde{p}$ of Example 3.4 is not a sum of bivariates	55
B.2. Pseudocode for TRW-S	57
C. Referenced Results	59

1. Introduction

In this work, we study the optimization of functions with $n > 2$ arguments that have a representation as a sum of several functions that have only 2 of the n arguments each. [Definition 1.1](#) introduces the rigorous definition that will be used throughout the work.

Definition 1.1 (Sum of Bivariates). *Let the function $F : \Omega \rightarrow \mathbb{R}$ be defined on a product space $\Omega := \Omega_1 \times \cdots \times \Omega_n$ of factors $\Omega_i \subset \mathbb{R}$ with a finite number of elements, where $|\Omega_i| = K_i$ for all $i \in \mathcal{V} := \mathbb{N}_{\leq n}$ and for some $n, K_i \in \mathbb{N}$.*

*We call the function F a **sum of bivariates** if there exists an index set $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$ and bivariate functions $f_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}, (i, j) \in \mathcal{E}$ with*

$$F(x_1, \dots, x_n) = \sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j) \quad \forall (x_1, \dots, x_n) \in \Omega.$$

The function model can be interpreted as indexed by a graph: The vertices index arguments and the edges index bivariate summands. If the functions $f_{i,j}, (i, j) \in \mathcal{E}$ were instead univariate, due to the additive structure, the problem would be separable into n individual univariate optimization problems. In this sense, the function model at hand presents the simplest non-trivial structure of its kind that makes no further assumptions on the functions $f_{i,j}, (i, j) \in \mathcal{E}$. Yet, it will turn out that both the structure of the bivariates $f_{i,j}, (i, j) \in \mathcal{E}$ as well as the edge set \mathcal{E} determine the hardness of optimizing F .

We restrict ourselves to finite candidate spaces Ω for the benefit of a simpler and tractable problem. Generalizations to uncountable candidate spaces are formally possible and practically useful.

The main focus of this work lies on so-called relaxation-based optimization methods for sums of bivariates. A *relaxation* is a transformation of a function $F : \Omega \rightarrow \mathbb{R}$ to a function

$$\mu \in \mathcal{M}_1^+(\Omega) \mapsto \int_{\Omega} F \, d\mu,$$

where (probability) measures $\mu \in \mathcal{M}_1^+(\Omega)$ replace the domain Ω and evaluate to the integral of F with respect to μ . Next to being linear functionals on a closed convex space and due to the sparse structure of sums of bivariates, relaxations of sums of bivariates will be shown to have useful properties for optimization. This perspective allows us to derive a wide range of optimization methods for sums of bivariates employing linear programming, coordinate ascent, dynamic programming, and even closed-form solution approaches. To this end, we also study the ℓ^2 -approximation and complexity theory of sums of bivariates.

Experiments on random functions, vertex coloring, and signal reconstruction problems provide insights into qualitatively different function classes that can be modeled as sums of bivariates.

Motivation. From a theoretical perspective, the motivation in studying the optimization of sums of bivariates lies in the pursuit of tractable models for global optimization. Making no assumptions on the structure of individual bivariates $f_{i,j}, (i, j) \in \mathcal{E}$ is an intriguing alternative to tractable models in the literature that often feature some sort of linear, convex, Lipschitz, or smooth structure.

From a practical perspective, sums of bivariates model a wide range of applications, where tractable instances often arise as inverse problems in signal reconstruction. In such problems, we often encounter a structure

$$F(x_1, \dots, x_n) = \sum_{i \in \mathbb{N}_{\leq n}} H_i(x_i) + \sum_{i,j \in \mathbb{N}_{\leq n}} G_{i,j}(x_i, x_j) \quad \forall (x_1, \dots, x_n) \in \Omega,$$

where $\Omega := \Omega_1 \times \cdots \times \Omega_n$ models the possible reconstructions and $F : \Omega \rightarrow \mathbb{R}$ models their quality based on argument-wise errors $H_i : \Omega_i \rightarrow \mathbb{R}, i \in \mathbb{N}_{\leq n}$ that are based on measured data, as well as pairwise regularization terms $G_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}, (i,j) \in \mathbb{N}_{\leq n}^2$. Clearly, such problems are sums of bivariates.

Sums of bivariates also appear in various other applications, such as, in the *Markowitz model for portfolio selection* as the hybrid objective that is the linear combination of portfolio risk and return [Mar52], and as the Hamiltonian of so-called *Sherrington–Kirkpatrick spin glasses* [SK75] as well as that of *Hopfield networks* [Ama72].

Due to its linearity and its commutativity with the projection operation for bivariate functions, relaxation is a particularly promising transformation for the function class at hand.

The pursuit of an elementary and consistent derivation of relaxation-based results from principles of mathematical optimization, approximation, and complexity theory inspires the focus of this work.

Related Work. Relaxation techniques for sums of bivariates were originally developed by [Sch76]. A review of associated linear programming formulations of optimization problems encountered in this work is given by [Wer07].

State-of-the-art methods that have a structure similar to those that we derive from relaxations of sums of bivariates have been studied in [Kol05; Kol14; Kap+15; Tou+18; Tou+20]. Subgradient methods that can serve as an alternative to the block coordinate ascent-based algorithms, which we introduce, have been developed by [SG07; KPT07]. Conversely, block-coordinate ascent methods, similar to those developed in this work, have also been used to create scalable solvers for integer linear programming [LS21]. Methods similar to our entropy-regularized relaxations have been developed before [GJ07; LI13].

Complexity results for various problem formulations encountered in this work are derived in [LSH16]. In particular, [PW15] show that any linear program can be reduced in linear time to what we call the *dual linear program* associated with a relaxation of sums of bivariates.

Extensive collections of results on the subject from a perspective of mathematical optimization as well as from a probabilistic one can be found in [Sav+19] and [WJ08], respectively. Very similar coordinate ascent-based optimization approaches for the Markov Random Field model have been described by [OP24].

Outline. In Section 2, we cover the basic ℓ^2 -approximation of functions by sums of bivariates. This leads us to derive an elementary characterization of parameterizations of zero and thereby the dual variables encountered later on.

We describe hard as well as non-trivial easy instances of sums of bivariates using complexity theory and dynamic programming in Section 3.1. An application of a no-free-lunch theorem will also be presented.

In Section 3.2, we focus on the central results on relaxations of sums of bivariates and the reconstructions of measures from bivariate marginals. We extend this analysis with a regularized version of relaxation for sums of bivariates in Section 3.3. We conclude our theory with a method that allows us to reconstruct solutions to the optimization of sums of bivariates from relaxation-associated problems in Section 3.4.

Section 4 contains the derivation of algorithms based on the various problem formulations and results encountered throughout the previous sections. Pseudocode for all discussed algorithms is included in this section.

The experiments in Section 5 consider the minimization of sums of random bivariates, vertex coloring, and signal reconstruction problems, which constitute qualitatively different problem classes. We verify qualitative properties of the algorithms and contrast hypotheses

- Given a product space $\Omega = \Omega_1 \times \cdots \times \Omega_n, n \in \mathbb{N}$, we define the projection $\pi_i^\Omega : \Omega \rightarrow \Omega_i, x \mapsto x_i$. Often we will drop the superscript and write $\pi_i := \pi_i^\Omega$, if the domain is clear by context. In the same setting, where additionally $\Omega_i = \mathbb{N}_{\leq m}$, we often denote $f \in \mathbb{R}^{\Omega_i}$ by $[f(1) \ \dots \ f(m)]_i$.
- A function $F : \Omega \rightarrow \mathbb{R}$ is called *separable* if there exist $f_i \in \mathbb{R}^{\Omega_i}, i \in \mathbb{N}_{\leq n}$ with $F \equiv f_1 + \cdots + f_n$.
- We often replace universally quantified variables by \cdot and do not explicitly denote the quantifier.
- Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with vertices $\mathcal{V} = \mathbb{N}_{\leq n}$ and edges $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$, we define $\mathcal{N}(i) := \{(i, j) \in \mathcal{E}\}$ and $\tilde{\mathcal{N}}(i) := \{j \in \mathcal{V} \mid (i, j) \in \mathcal{E} \vee (j, i) \in \mathcal{E}\}, i \in \mathcal{V}$. In oriented trees, we will assume that a parent vertex is in the left factor of an edge and a child vertex to be in the right factor of an edge.
- For a set Ω with $|\Omega| \in \mathbb{N}$ and $f \in \mathbb{R}^\Omega$, we define as $\text{lse}_\varepsilon(f) := \varepsilon \log \sum_{x \in \Omega} \exp(f(x)/\varepsilon)$ the ε -LogSumExp for all $\varepsilon > 0$. We will also use the notation $\text{lse}_\varepsilon^x(f(x)) := \text{lse}_\varepsilon(f)$ if f has further parameters.

2. Approximation

Initially, we would like to better understand the parameterization that is implicit to the defining functional equation of sums of bivariate in [Definition 1.1](#). This will yield not only a characterization of the defining functional equation but also a way to efficiently approximate functions by sums of bivariate, i.e. to project onto the space of sums of bivariate.

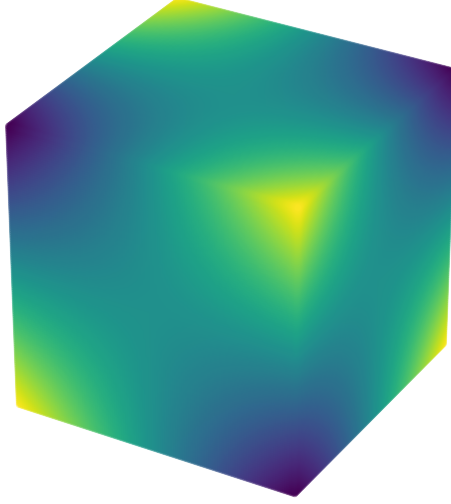


Figure 2: Visualization of a function L^2 -orthogonal to the sums of bivariate on $\Omega = [0, 1]^3$.

For a given function $G : \mathbb{N}_{\leq m}^n \rightarrow \mathbb{R}$, we call the function that sums G over all values that all but its k -th and ℓ -th argument can take, where $k < \ell$, the (k, ℓ) -marginal of G , i.e.

$$a, b \in \mathbb{N}_{\leq m} \mapsto \sum_{\substack{y_i \in \mathbb{N}_{\leq m} \\ i \in \mathbb{N}_{\leq n} \\ i \neq k, \ell}} G(y_1, \dots, y_{k-1}, a, y_{k+1}, \dots, y_{\ell-1}, b, y_{\ell+1}, \dots, y_n).$$

The first theorem provides a necessary and sufficient condition for the approximation of a function G by sums of bivariate in terms of all (k, ℓ) -marginals. By that, the theorem also provides a characterization of the defining functional when replacing G by a function that is postulated to be a sum of bivariate. The theorem features a more compact notation of the (k, ℓ) -marginal.

Insights into conditions and limits for bivariate to be the marginals of a possibly unknown function G will be covered in [Theorem 3.4](#) and [Example 3.5](#).

Theorem 2.1 (Approximation). *Given $G : \mathbb{N}_{\leq m}^n \rightarrow \mathbb{R}$, $m \in \mathbb{N}$, a sum of bivariate $F : \mathbb{N}_{\leq m}^n \rightarrow \mathbb{R}$, where $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$, $\mathcal{E} = \{(i,j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$, and $\mathcal{V} = \mathbb{N}_{\leq n}$, is a best ℓ^2 -approximation of G if and only if*

$$\sum_{y \in \mathbb{N}_{\leq m}^{n-2}} G(z^{k,\ell}(a, b, y)) = \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} F(z^{k,\ell}(a, b, y)) \quad \forall a, b \in \mathbb{N}_{\leq m} \quad \forall (k, \ell) \in \mathcal{E}$$

$$\begin{aligned}
&= m^{n-2} f_{k,\ell}(a, b) + m^{n-3} \left(\left(\sum_{j < k} \overline{f_{j,k}}^1(a) + \overline{f_{j,\ell}}^1(b) \right) \right. \\
&\quad \left. + \left(\sum_{k < j < \ell} \overline{f_{k,j}}^2(a) + \overline{f_{j,\ell}}^1(b) \right) \right. \\
&\quad \left. + \left(\sum_{\ell < j} \overline{f_{k,j}}^2(a) + \overline{f_{\ell,j}}^2(b) \right) \right) \\
&\quad + m^{n-4} \sum_{i, j \notin \{k, \ell\}} \overline{\overline{f_{i,j}}}, \quad (\text{this eq. only holds if } n \geq 4)
\end{aligned}$$

where

$$\begin{aligned}
z^{k,\ell}(a, b, y) &:= (y_1, \dots, y_{k-1}, a, y_k, \dots, y_{\ell-2}, b, y_{\ell-1}, \dots, y_{n-2}), \\
\overline{f_{i,j}}^1(s) &:= \sum_{x \in \mathbb{N}_{\leq m}} f_{i,j}(x, s), \quad \overline{f_{i,j}}^2(s) := \sum_{x \in \mathbb{N}_{\leq m}} f_{i,j}(s, x), \quad \forall s \in \{a, b\}, \text{ and} \\
\overline{\overline{f_{i,j}}} &:= \sum_{x, s \in \mathbb{N}_{\leq m}} f_{i,j}(x, s), \quad \forall (i, j) \in \mathcal{E}.
\end{aligned}$$

Proof. We want to find a bivariate function $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$, where $\mathcal{E} = \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$, and $\mathcal{V} = \mathbb{N}_{\leq m}$ that minimizes

$$\|G - F\|_{\ell^2} = \sqrt{\sum_{x \in \mathbb{N}_{\leq m}^n} \left(G(x) - \sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j) \right)^2}.$$

Due to convexity, the first-order condition

$$\begin{aligned}
0 &\stackrel{!}{=} \frac{\partial}{\partial f_{k,\ell}(a, b)} \sum_{x \in \mathbb{N}_{\leq m}^n} \left(G(x) - \sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j) \right)^2 \quad \forall a, b \in \mathbb{N}_{\leq m} \quad \forall (k, \ell) \in \mathcal{E} \\
&\quad \text{(first-order condition)} \\
&= \frac{\partial}{\partial f_{k,\ell}(a, b)} \sum_{y_k, y_\ell \in \mathbb{N}_{\leq m}} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \left(G(z^{k,\ell}(y_k, y_\ell, y)) - \sum_{(i,j) \in \mathcal{E}} f_{i,j}(z_i^{k,\ell}(y_k, y_\ell, y), z_j^{k,\ell}(y_k, y_\ell, y)) \right)^2 \\
&\quad \text{(reorder summation; definition } z^{k,\ell}) \\
&= -2 \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} G(z^{k,\ell}(a, b, y)) - \sum_{(i,j) \in \mathcal{E}} f_{i,j}(z_i^{k,\ell}(a, b, y), z_j^{k,\ell}(a, b, y)) \quad \text{(differentiation)} \\
&\iff \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} G(z^{k,\ell}(a, b, y)) = \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{(i,j) \in \mathcal{E}} f_{i,j}(z_i^{k,\ell}(a, b, y), z_j^{k,\ell}(a, b, y)) \\
&\quad \text{(multiply by } -1/2; \text{ add marginalized sum of bivariates)} \\
&= \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} F(z^{k,\ell}(a, b, y)) \quad \text{(definition of } F) \\
&= m^{n-2} f_{k,\ell}(a, b) + m^{n-3} \left(\left(\sum_{j < k} \overline{f_{j,k}}^1(a) + \overline{f_{j,\ell}}^1(b) \right) \right. \\
&\quad \left. + \left(\sum_{k < j < \ell} \overline{f_{k,j}}^2(a) + \overline{f_{j,\ell}}^1(b) \right) \right. \\
&\quad \left. + \left(\sum_{\ell < j} \overline{f_{k,j}}^2(a) + \overline{f_{\ell,j}}^2(b) \right) \right)
\end{aligned}$$

$$+ m^{n-4} \sum_{i,j \notin \{k,\ell\}} \overline{f_{i,j}}$$

(reorder summation based on occurrences of a and b ; definitions $\overline{f_{\cdot,\cdot}}$, $\overline{f_{\cdot,\cdot}}$)

is necessary and sufficient. \square

Figure 2 shows the residual of (numerically) approximating $(x, y, z) \in [0, 1]^3 \mapsto xyz$ by sums of bivariates. The resulting function has a ℓ^2 -best approximant by sums of bivariates that is the function 0, therefore, it is ℓ^2 -orthogonal to the sums of bivariates. There are, however, still multiple sets of bivariates that sum to 0 and, therefore, parameterize the best approximant.

Thus, Corollary 2.1 characterizes the parameterizations of 0 by sums of bivariates. As we will see in Sections 3 and 4, the result is central to solvers that use dual formulations and—sometimes without loss—transforms the search space of the minimization of sums of bivariates into one with comparatively low dimensions.

A relevant feature of the following result is the representation not only of the constant sum of bivariates F as a sum of univariates, but, in particular, a representation of the individual bivariates $f_{i,j}$ as a sum of univariates.

Corollary 2.1 (Dual Variables). *In the setting of Theorem 2.1, we have*

$$\begin{aligned} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} F(z^{k,\ell}(a, b, y)) &= 0 \quad \forall a, b \in \{1, \dots, m\} \quad \forall (k, \ell) \in \mathcal{E} \quad (\text{marginal null constraint}) \\ \iff F &\equiv 0 \quad (\text{global null constraint}) \\ \iff \forall (i, j) \in \mathcal{E} : \exists \rho_{i,j}, \rho_{j,i} \in \mathbb{R}^{\mathbb{N}_{\leq m}} : &\left(f_{i,j} \equiv \rho_{i,j}(\cdot_i) + \rho_{j,i}(\cdot_j) \right) \\ &\wedge \left(\forall i \in \mathcal{V} : \exists \rho_i \in \mathbb{R} : \sum_{j \in \mathcal{V}} \rho_{i,j} \equiv \rho_i \right) \\ &\wedge \left(\sum_{i \in \mathcal{V}} \rho_i = 0 \right). \quad (\text{dual null constraint}) \end{aligned}$$

Remark 1. Clearly, a similar result holds for constant sums of bivariates.

Proof. We prove the statement as a Ringschluss.

“marginal null constraint \Rightarrow dual null constraint”: Consider, that for all $(k, \ell) \in \mathcal{E}$ and for all $a, b \in \mathbb{N}_{\leq m}$, we have

$$\begin{aligned} 0 &= \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} F(z^{k,\ell}(a, b, y)) \quad (\text{marginal null constraint}) \\ &= \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{(i,j) \in \mathcal{E}} f_{i,j}(z_i^{k,\ell}(a, b, y), z_j^{k,\ell}(a, b, y)) \\ &\quad (\text{representation of sum of bivariates}) \\ \iff f_{k,\ell}(a, b) &= -\frac{1}{m^{n-2}} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{(i,j) \in \mathcal{E} \setminus \{(k,\ell)\}} f_{i,j}(z_i^{k,\ell}(a, b, y), z_j^{k,\ell}(a, b, y)) \\ &\quad (\text{multiplication by } -1/m^{n-2}; \text{ addition of } f_{k,\ell}(a, b)) \end{aligned}$$

Since the right-hand side of the previous equation is a sum of univariate functions in either a or b , or of constants, for all $k, \ell \in \mathcal{E}$, we get that

$$\exists \rho_{k,\ell}, \rho_{\ell,k} \in \mathbb{R}^{\mathbb{N}_{\leq m}} : f_{k,\ell} \equiv \rho_{k,\ell} + \rho_{\ell,k}.$$

Further, for all $a \in \mathbb{N}_{\leq n}$ and for all $k \in \mathcal{V}$, we have

$$\begin{aligned} 0 &= \sum_{b \in \mathbb{N}_{\leq m}} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} F(z^{k,l}(a, b, y)) && \text{(sum of marginal null constraint)} \\ &= \sum_{b \in \mathbb{N}_{\leq m}} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{(i,j) \in \mathcal{E}} f_{i,j}(z_i^{k,\ell}(a, b, y), z_j^{k,\ell}(a, b, y)) && \text{(representation of sum of bivariates)} \\ &= \sum_{b \in \mathbb{N}_{\leq m}} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{(i,j) \in \mathcal{E}} \rho_{i,j}(z_i^{k,\ell}(a, b, y)) + \rho_{j,i}(z_j^{k,\ell}(a, b, y)) && \text{(previous result)} \\ &= \sum_{b \in \mathbb{N}_{\leq m}} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} \rho_{i,j}(z_i^{k,\ell}(a, b, y)) && \text{(reorder summation)} \\ \iff \sum_{j \in \mathcal{V}} \rho_{k,j}(a) &= -\frac{1}{m^{n-1}} \sum_{b \in \mathbb{N}_{\leq m}} \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{i \in \mathcal{V} \setminus \{k\}} \sum_{j \in \mathcal{V}} \rho_{i,j}(z_i^{k,\ell}(a, b, y)). && \text{(multiplication by } -1/m^{n-1}; \text{ addition of } \sum_{j \in \mathcal{V}} \rho_{k,j}(a)) \end{aligned}$$

Since the right-hand side of the previous equation is a constant, and by a similar argument for the variable b , we get for all $k \in \mathcal{V}$ that

$$\exists \rho_k \in \mathbb{R} : \sum_{j \in \mathcal{V}} \rho_{k,j} \equiv \rho_k.$$

Finally, we get the last claim for an arbitrary choice of $a, b \in \mathbb{N}_{\leq m}$ and $k, \ell \in \mathcal{E}$ by

$$\begin{aligned} 0 &= \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} F(z^{k,l}(a, b, y)) && \text{(marginal null constraint)} \\ &= \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{(i,j) \in \mathcal{E}} \rho_{i,j}(z_i^{k,\ell}(a, b, y)) + \rho_{j,i}(z_j^{k,\ell}(a, b, y)) && \text{(previous result)} \\ &= \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{i \in \mathcal{V}} \left(\sum_{j \in \mathcal{V}} \rho_{i,j}(z_i^{k,\ell}(a, b, y)) \right) && \text{(reorder summation)} \\ &\equiv \sum_{y \in \mathbb{N}_{\leq m}^{n-2}} \sum_{i \in \mathcal{V}} \rho_i && \text{(previous result)} \\ &= \cancel{m^{n-2}} \sum_{i \in \mathcal{V}} \rho_i. && \text{(constant summation)} \end{aligned}$$

“dual null constraint \Rightarrow global null constraint”: We have for all $x_1, \dots, x_n \in \mathbb{N}_{\leq m}$ that

$$\begin{aligned}
 F(x_1, \dots, x_n) &= \sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j) && \text{(representation of sum of bivariates)} \\
 &= \sum_{(i,j) \in \mathcal{E}} \rho_{i,j}(x_i) + \rho_{j,i}(x_j) && \text{(first property of dual null constraint)} \\
 &= \sum_{i \in \mathcal{V}} \left(\sum_{j \in \mathcal{V}} \rho_{i,j}(x_i) \right) && \text{(reorder summation)} \\
 &\equiv \sum_{i \in \mathcal{V}} \rho_i && \text{(second property of dual null constraint)} \\
 &= 0. && \text{(third property of dual null constraint)}
 \end{aligned}$$

“global null constraint \Rightarrow marginal null constraint”: Clear. \square

Since the main goal of this work is the analysis of relaxation-based solvers, further approximation-based insights into sums of bivariates will be deferred.

3. Optimization

3.1. Fundamental Perspective

We now attempt to classify as well as understand some hard- and easy-to-optimize instances of sums of bivariates. Such insight may help in modeling, managing our expectations for the optimization methods we design, and in interpreting their performance.

The Hamiltonian cycle problem asks to determine whether for a given graph there exists a sequence of successively adjacent vertices that are nonrepetitive. The Hamiltonian cycle problem is among the hardest decision problems for which a solution is polynomially sized and efficiently verifiable.

In [Example 3.1](#), we recite the reduction of the Hamiltonian cycle problem to checking whether the minimal value of a sum of bivariates is 0. By that we show that efficiently minimizing sums of bivariates is at least as hard as such hardest decision problems.

Example 3.1 (The Hamiltonian Cycle Problem as a Sum of Bivariates [[Sav+19](#), p. 13 Sec. 1.3]).

Let $(\mathcal{V}', \mathcal{E}')$ be a graph for which we want prove existence of a Hamiltonian cycle.

Define $\mathcal{V} := \mathbb{N}_{\leq n}$, $n = |\mathcal{V}'|$, $\mathcal{E} = \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$, $\Omega_i := \mathcal{V}'$, $i \in \mathbb{N}_{\leq n}$, $\Omega := \Omega_1 \times \cdots \times \Omega_n$, and for all $(i, j) \in \mathcal{E}$ for all $(x_i, x_j) \in \Omega_i \times \Omega_j$, that

$$f_{i,j}(x_i, x_j) := \begin{cases} 0 & \text{if } x_i \neq x_j \wedge (i+1 = j \implies (x_i, x_j) \in \mathcal{E}') \\ 1 & \text{else.} \end{cases}$$

Therefore,

$$\begin{aligned} \min_{x \in \Omega} \sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j) &= 0 \\ \iff \exists x \in \Omega : &\left(\forall (i, j) \in \mathcal{E} : x_i \neq x_j \wedge (i+1 = j \implies (x_i, x_j) \in \mathcal{E}') \right) \\ \iff &(\mathcal{V}', \mathcal{E}') \text{ has a Hamiltonian cycle.} \end{aligned}$$

The problem of deciding whether the graph $\mathcal{G}' = (\mathcal{V}', \mathcal{E}')$ has a Hamiltonian cycle is NP-complete. Therefore, since the Hamiltonian cycle problem can be reduced in polynomial time to checking the minimal value of a sum of bivariates, we know that minimization of the sum of bivariates is NP-hard.

Next, we show in [Example 3.2](#) that the minimization of sums of bivariates can be efficiently reduced to *integer linear programming*. This implies that minimizing rationally-valued sums of bivariates is also NP-easy. In conjunction with the result of NP-hardness by [Example 3.1](#), we therefore know that the minimization of rationally-valued sums of bivariates is NP-equivalent. This means that there exists an exact and efficient algorithm for the minimization of rationally-valued sums of bivariates if and only if there is an efficient algorithm for the hardest decision problems for which solutions are polynomially sized and efficiently verifiable, i.e. if $P = NP$ [[KV06](#), p. 368 Prop. 15.35].

Example 3.2 (Sums of Bivariates as Integer Linear Programs). In the setting of [Definition 1.1](#), let $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j} : \Omega_1 \times \cdots \times \Omega_n \rightarrow \mathbb{Q}$ encode a sum of bivariates.

Let, w.l.o.g., $|\Omega_1| = \cdots = |\Omega_n| = m \in \mathbb{N}$. Consider the binary variables of the integer linear program to be

$$x_{i,j} \in \{0, 1\}^{m \times m} \quad \forall (i, j) \in \mathcal{E},$$

and the constraints to be

$$\begin{cases} \sum_{s,t=1}^m x_{i,j;s,t} = 1 & \forall (i,j) \in \mathcal{E} \\ \sum_{s=1}^m x_{i,j;s,t} = \sum_{s=1}^m x_{k,j;s,t} & \forall t \in \mathbb{N}_{\leq m} \forall (i,j), (k,j) \in \mathcal{E} \\ \sum_{s=1}^m x_{i,j;t,s} = \sum_{s=1}^m x_{k,j;t,s} & \forall t \in \mathbb{N}_{\leq m} \forall (i,j), (i,k) \in \mathcal{E} \\ \sum_{s=1}^m x_{i,j;s,t} = \sum_{s=1}^m x_{j,k;t,s} & \forall t \in \mathbb{N}_{\leq m} \forall (i,j), (j,k) \in \mathcal{E}. \end{cases}$$

The objective to be minimized can be modeled as

$$\sum_{(i,j) \in \mathcal{E}} \sum_{s,t=1}^m f_{i,j}(s,t) \cdot x_{i,j;s,t},$$

which leaves us with a linear integer program that has a polynomial size w.r.t. to the size of the sum of bivariates. Clearly, we obtain an optimal solution $y^* \in \Omega$ to the minimization of sums of bivariates given an optimal solution to the integer linear program $x^* \in (\{0,1\}^{m \times m})^{|\mathcal{E}|}$ by picking

$$y_i^* \in \arg \max_{s \in \mathbb{N}_{\leq m}} \sum_{t=1}^m x_{i,j;s,t}^* \quad \text{and} \quad y_j^* \in \arg \max_{j \in \mathbb{N}_{\leq m}} \sum_{s=1}^m x_{i,j;s,t}^* \quad \forall (i,j) \in \mathcal{E},$$

which is well-defined by our constraints.

In contrast to the described hard instances, we now introduce relatively easy-to-minimize instances of sums of bivariates, which are widely known. Interestingly, this class of functions is not characterized by requirements on the structure of the bivariates but on the structure that encodes the dependencies, that is, the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ that indexes the sum of bivariates $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$ in Definition 1.1. Such easy instances are, in fact, characterized by tree-structured index graphs \mathcal{G} .

The following Theorem 3.1 is constructive in the sense that it describes an algorithm to efficiently solve such easy instances. The minimization of sums of bivariates with these tree-structured indices can be interpreted as a dynamic program that sequentially allows the reduction of the problem to a smaller instance.

Theorem 3.1 (Dynamic Programming for Tree-Indexed Instances). *Let $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$ be a sum of bivariates, where $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a tree, and let $(i_t, j_t) \in \mathcal{E}, j_t := |\mathcal{V}| + 1 - t, t \in \mathbb{N}_{\leq |\mathcal{V}| - 1}$ be an edge counting, such that*

- $\mathcal{E}_t := \mathcal{E}_{t-1} \setminus \{(i_t, j_t)\}$, where $\mathcal{E}_0 := \mathcal{E} \cup \{(*, i_{|\mathcal{V}| - 1})\}$ and $*$ is a placeholder symbol,
- $\mathcal{V}_t := \mathcal{V}_{t-1} \setminus \{j_t\}$, where $\mathcal{V}_0 := \mathcal{V}$,
- j_t is a leaf of $\mathcal{G}_{t-1} := (\mathcal{V}_{t-1}, \mathcal{E}_{t-1})$,
- define $f_{0;*, i_{|\mathcal{V}| - 1}} := 0$, $f_{0;k,\ell} := f_{k,\ell}$ and

$$f_{t;k,\ell} := \begin{cases} f_{t-1;k,\ell} + \min_{x_{j_t}} f_{t-1;i_t,j_t}(\cdot, x_{j_t}) & \text{if } \ell = i_t \\ f_{t-1;k,\ell} & \text{else,} \end{cases}$$

for all $(k, \ell) \in \mathcal{E}_t$ for all $t \in \mathbb{N}_{\leq |\mathcal{V}| - 1}$, and

- $F_t := \sum_{(i,j) \in \mathcal{E}_t} f_{t;i,j}$ for all $t \in \mathbb{N}_{[0, |\mathcal{V}| - 1]}$.

Then,

- $\min F = \min F_t$ for all $t \in \mathbb{N}_{\leq |\mathcal{V}|-1}$, and
- $x^* \in \arg \min_{x \in \Omega} F(x)$ if $x_t^* \in \arg \min_{x_t \in \Omega_t} F_{|\mathcal{V}|-t}(x_1^*, \dots, x_{t-1}^*, x_t)$ for all $t \in \mathbb{N}_{\leq |\mathcal{V}|}$.

Proof. The first result can be inductively concluded from the following equation. We have for all $t \in \mathbb{N}_{\leq |\mathcal{V}|-1}$ that

$$\begin{aligned}
F_t &= \sum_{(i,j) \in \mathcal{E}_t} f_{t;i,j} && \text{(definition of } F_t) \\
&= \left(\sum_{(i,j) \in \mathcal{E}_t} f_{t-1;i,j} \right) + \sum_{(i,j) \in \mathcal{E}_t} f_{t;i,j} - f_{t-1;i,j} && \text{(reorder summation; partition of zero)} \\
&= \left(\sum_{(i,j) \in \mathcal{E}_t} f_{t-1;i,j} \right) + \min_{x_{j_t}} f_{t-1;i_t,j_t}(\cdot, x_{j_t}) && \text{(definition of } f_{t;\cdot,\cdot}) \\
&= \min_{x_{j_t}} \sum_{(i,j) \in \mathcal{E}_{t-1}} f_{t-1;i,j} && \text{(definition of } \mathcal{E}_t; j_t \text{ is leaf of } \mathcal{G}_{t-1}) \\
&= \min_{x_{j_t}} F_{t-1}. && \text{(definition of } F_{t-1})
\end{aligned}$$

The second result is implied by

$$(x_1^*, \dots, x_t^*) \in \arg \min_{x_1, \dots, x_t} F_{|\mathcal{V}|-t}(x_1^*, \dots, x_{t-1}^*, x_t) \quad \forall t \in \mathbb{N}_{\leq |\mathcal{V}|},$$

which is also true inductively by

$$\begin{aligned}
F_{|\mathcal{V}|-t}(x_1^*, \dots, x_t^*) &= \min_{x_t} F_{|\mathcal{V}|-t}(x_1^*, \dots, x_{t-1}^*, x_t) && \text{(definition of } x_t^*) \\
&= F_{|\mathcal{V}|-t}(x_1^*, \dots, x_{t-1}^*) && \text{(previous equation)} \\
&= \min F_{|\mathcal{V}|-t} && \text{(induction hypothesis)} \\
&= \min F_{|\mathcal{V}|-t}, && \text{(first result)}
\end{aligned}$$

where $t \in \mathbb{N}_{[2, |\mathcal{V}|]}$ and $x_1^* \in \arg \min_{x_1} F_{|\mathcal{V}|-1}(x_1)$ by definition. \square

In [Example 3.3](#), we present a small but detailed example of a sum of bivariate functions that we minimize using the method described in [Theorem 3.1](#).

Example 3.3 (Dynamic Programming for Tree-Indexed Instances). *In the setting of [Definition 1.1](#), consider*

$$\begin{aligned}
\mathcal{V} &:= \{1, 2, 3, 4, 5\}, \quad \mathcal{E} := \{(1, 2), (2, 3), (3, 4), (3, 5)\}, \\
\Omega_1 &= \dots = \Omega_5 = \{1, 2, 3\},
\end{aligned}$$

and a sum of bivariate functions

$$F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j},$$

where

$$\begin{aligned}
f_{1,2} &:= \begin{bmatrix} 3 & 8 & 0 \\ 6 & 8 & 9 \\ 6 & 1 & 9 \end{bmatrix} \in \mathbb{R}^{\Omega_1 \times \Omega_2}, & f_{2,3} &:= \begin{bmatrix} 8 & 4 & 2 \\ 6 & 9 & 9 \\ 6 & 4 & 2 \end{bmatrix} \in \mathbb{R}^{\Omega_2 \times \Omega_3}, \\
f_{3,4} &:= \begin{bmatrix} 7 & 0 & 0 \\ 2 & 2 & 1 \\ 4 & 2 & 6 \end{bmatrix} \in \mathbb{R}^{\Omega_3 \times \Omega_4}, & f_{3,5} &:= \begin{bmatrix} 8 & 9 & 6 \\ 0 & 2 & 7 \\ 9 & 3 & 3 \end{bmatrix} \in \mathbb{R}^{\Omega_3 \times \Omega_5}.
\end{aligned}$$

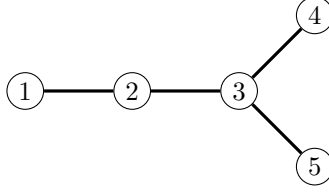


Figure 3: Visualization of the variable dependencies of the sum of bivariate of Example 3.3.

Then, in the setting of Theorem 3.1, we have

$$j_1 = 5, \quad j_2 = 4, \quad j_3 = 3, \quad \text{and} \quad j_4 = 2,$$

as well as

$$\begin{aligned}
F_0 &\equiv f_{1,2} + f_{2,3} + f_{3,4} + f_{3,5} \\
F_1 &\equiv f_{1,2} + f_{2,3} + f_{3,4} + \min_{x_5 \in \Omega_5} f_{3,5}(\cdot, x_5) \\
&\equiv f_{1,2} + f_{2,3} + f_{3,4} + [6 \quad 0 \quad 3]_3 \\
F_2 &\equiv f_{1,2} + f_{2,3} + \left(\min_{x_4 \in \Omega_4} f_{3,4}(\cdot, x_4) \right) + \left(\min_{x_5 \in \Omega_5} f_{3,5}(\cdot, x_5) \right) \\
&\equiv f_{1,2} + f_{2,3} + [0 \quad 1 \quad 2]_3 + [6 \quad 0 \quad 3]_3 \\
F_3 &\equiv f_{1,2} + \left(\min_{x_3 \in \Omega_3} f_{2,3}(\cdot, x_3) + \left(\min_{x_4 \in \Omega_4} f_{3,4}(x_3, x_4) \right) + \left(\min_{x_5 \in \Omega_5} f_{3,5}(x_3, x_5) \right) \right) \\
&\equiv f_{1,2} + \left(\min_{x_3 \in \Omega_3} f_{2,3}(\cdot, x_3) + [6 \quad 1 \quad 5]_3(x_3) \right) \\
&\equiv f_{1,2} + \min_{x_3 \in \Omega_3} \begin{bmatrix} 14 & 5 & 7 \\ 12 & 10 & 14 \\ 12 & 5 & 7 \end{bmatrix}_{2,3}(\cdot, x_3) \\
&\equiv f_{1,2} + [5 \quad 10 \quad 5]_2 \\
F_4 &\equiv \min_{x_2 \in \Omega_2} f_{1,2}(\cdot, x_2) + \left(\min_{x_3 \in \Omega_3} f_{2,3}(x_2, x_3) + \left(\min_{x_4 \in \Omega_4} f_{3,4}(x_3, x_4) \right) + \left(\min_{x_5 \in \Omega_5} f_{3,5}(x_3, x_5) \right) \right) \\
&\equiv \min_{x_2 \in \Omega_2} \begin{bmatrix} 8 & 18 & 5 \\ 11 & 18 & 14 \\ 11 & 11 & 14 \end{bmatrix}_{1,2}(\cdot, x_2) \\
&\equiv [5 \quad 11 \quad 11]_1.
\end{aligned}$$

By Theorem 3.1, we know $\min F = \min F_4 = 5$ and for

$$\begin{aligned}
x_1^* &\in \arg \min_{x_1 \in \Omega_1} F_4(x_1) = \{1\}, \\
x_2^* &\in \arg \min_{x_2 \in \Omega_2} F_3(x_1^*, x_2) = \arg \min_{x_2 \in \Omega_2} [8 \quad 18 \quad 5]_2 = \{3\}, \\
x_3^* &\in \arg \min_{x_3 \in \Omega_3} F_2(x_1^*, x_2^*, x_3) = \arg \min_{x_3 \in \Omega_3} [12 \quad 5 \quad 7]_3 = \{2\}, \\
x_4^* &\in \arg \min_{x_4 \in \Omega_4} F_1(x_1^*, x_2^*, x_3^*, x_4) = \arg \min_{x_4 \in \Omega_4} [6 \quad 6 \quad 5]_4 = \{3\}, \text{ and} \\
x_5^* &\in \arg \min_{x_5 \in \Omega_5} F(x_1^*, x_2^*, x_3^*, x_4^*, x_5) = \arg \min_{x_5 \in \Omega_5} [6 \quad 8 \quad 13]_5 = \{1\},
\end{aligned}$$

we have $(x_1^*, \dots, x_5^*) \in \arg \min F$.

We identified some hard and easy instances of sums of bivariate functions. However, it does not have to be our goal to solve the hardest problems optimally. We may be happy to solve other instances overproportionally efficient with a particular optimization algorithm.

Related to this perspective are results that constrain the overperformance of some algorithms over others when applied to a class of objective functions, which are called *no-free-lunch results*. Such a result could, in principle, give us a formal reason to abandon the goal of minimizing general sums of bivariate functions—or at the very least a reason to constrain to particular subclass of sums of bivariate functions. The following [Theorem 3.2](#), presented rather informally, gives us a characterization of classes of objectives for which all arguably reasonable algorithms have the same performance distribution with respect to all arguably reasonable performance measures.

Theorem 3.2 ((No-)Free-Lunch Theorem [IT04, p. 3 Thm. 2] & [SVW01]). *Let $\Omega, Y \subset \mathbb{R}$ be finite sets. If and only if $\mathcal{F} \subseteq Y^\Omega$ is closed under composition with permutations of Ω , then for*

- *any two algorithms a, b , mapping (candidate, objective value)-sequences to non-repeating successive candidate-objective value points,*
- *any budget $k \in \mathbb{N}_{\leq |\Omega|}$, and*
- *any performance measure, mapping objective value-sequences to real numbers,*

the distribution of performances of length- k sequences', generated by a and b , when applied to \mathcal{F} , is equivalent.

We conclude this section with the construction of a counterexample. That is, we show in [Example 3.4](#) that sums of bivariate functions are not closed under permutation of their domain. This means that we find an instance of the sums of bivariate functions $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j} : \Omega \rightarrow \mathbb{R}$, such that when we compose it with a permutation of Ω , we obtain a function that does not have a representation as a sum of bivariate functions.

Admittedly, the negation of [Theorem 3.2](#) is rather weak: It just implies that a performance measure exists, such that there exist two algorithms with non-equal performance measures when applied to sums of bivariate functions. However, this fact may still serve as a sanity check, as we would hope this is true if the endeavor of finding good optimizers for arbitrary sums of bivariate functions is meant to be worthwhile.

Example 3.4 (Sums of Bivariate Functions Not Closed Under Permutation). *Consider*

$$\begin{cases} F : \{0, 1\}^6 \longrightarrow \mathbb{Z}_{[-64, 64]} \\ x_1 x_2 \dots x_6 \longmapsto \sum_{\substack{i, j \in \mathbb{N}_{\leq 6} \\ i < j}} \chi_{x_i = x_j = 1}(x_i, x_j) \end{cases}$$

and the permutation of $\mathbb{N}_{[0, 63]}$ in cycle notation given by

$$p = (0, 8, 18, 11, 54, 41, 28, 26, 55, 59, 48, 40, 60, 24, 47, 12, 33, 63, 13, 22, 25, 16, 23, 32, 7, 36, 21, \\ 6, 1, 52, 44, 50, 42, 17, 10, 53, 37, 14, 39, 9, 58, 46, 38, 51, 5, 27, 56, 31, 15, 49, 35, 61, 45, 3, \\ 30, 19, 57, 34, 4, 43, 2, 62, 20, 29).$$

Then, let $\iota : \{0, 1\}^6 \rightarrow \mathbb{N}_{[0, 63]}$ be the decimal coding, then $\tilde{p} = \iota^{-1} \circ p \circ \iota$ is a permutation of $\{0, 1\}^6$, and $F \circ \tilde{p}$ not a sum of bivariate functions. The claim can be verified by running the program in [Appendix B.1](#).

3.2. Relaxation Principle

The focus of this work are relaxation-based optimization algorithms. Relaxation is a transformation of an optimization problem $\min_{x \in \Omega} F(x)$ that generalizes the concept of a candidate $x \in \Omega$ to probability measures $\mathcal{M}_1^+(\Omega)$, as well as the objective value of a candidate from $F(x)$ to

$$\langle F, \mu \rangle := \int_{\Omega} F \, d\mu,$$

and is applicable to integrable objectives. Relaxations are used in other domains of (applied) mathematics to generalize solution concepts, such as in game theory and PDE-analysis. Although the transformation is often not of practical interest, it associates the original optimization problem with a unique linear optimization problem on a convex and often compact domain.

For our purposes, relaxations serve us in generating insight into our problem as well as to derive practically useful problem formulations.

The following [Theorem 3.3](#) recalls that one can generally associate a relaxed optimization problem with the problem of minimizing an objective function, such that it has the same minimal value and any optimal measure evaluates to zero on the set of suboptimal candidates of the original problem. This is, in particular, true for sums of bivariates.

The deep insight of [Theorem 3.3](#) is that given a candidate $\mu \in \mathcal{M}_1^+(\Omega)$ of the relaxed (global) problem, it is sufficient to know for all $(i, j) \in \mathcal{E}$ its bivariate marginal measure

$$\begin{aligned} \mu_{i,j} &:= \mu(\Omega_1 \times \cdots \times \Omega_{i-1} \times (\cdot) \times \Omega_{i+1} \times \cdots \times \Omega_{j-1} \times (\cdot) \times \Omega_{j+1} \times \cdots \times \Omega_n) \\ &= \mu \circ \pi_{i,j}^{-1} \in \mathcal{M}_1^+(\Omega_i \times \Omega_j), \end{aligned}$$

to determine its objective value. Here, the function $\pi_{i,j} : \Omega \rightarrow \Omega_i \times \Omega_j$ denotes the projection to the coordinates i, j . However, only if the graph that indexes the sum of bivariates is a tree, we associate optimal solutions of the relaxation with optimal solutions of the minimization of the sums of bivariates based solely on their bivariate marginal measures. By this we obtain a formulation for the optimization of tree-structured sums of bivariates that is already efficiently solvable and is an alternative to the dynamic programming approach presented in [Theorem 3.1](#).

Additionally, in case the sum of bivariates is indexed by a star graph, one obtains an even simpler linear programming formulation, which we will later learn to admit a closed-form solution.

Theorem 3.3 (Relaxation). *If a sum of bivariates $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j} : \Omega \rightarrow \mathbb{R}$ has a unique minimum, then*

$$\begin{aligned} &(\arg) \min_{x \in \Omega} F(x_1, \dots, x_n) \\ &= (\arg) \min_{\mu \in \mathcal{M}_1^+(\Omega)} \langle F, \mu \rangle. \end{aligned} \tag{global relaxation}$$

Further, if $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a tree, then

$$= \begin{cases} (\arg) \min_{\substack{\mu_{\cdot,\cdot} \in \mathcal{M}^+(\Omega \times \Omega) \\ \mu_{\cdot,\cdot} \in \mathcal{M}_1^+(\Omega)}} \sum_{(i,j) \in \mathcal{E}} \langle f_{i,j}, \mu_{i,j} \rangle \\ \text{s.t.} \quad \mu_i = \mu_{i,j} \circ \pi_i^{-1}, \quad \forall (i,j) \in \mathcal{E} \\ \quad \mu_j = \mu_{i,j} \circ \pi_j^{-1}, \quad \forall (i,j) \in \mathcal{E}. \end{cases} \tag{tree relaxation}$$

Further, if $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a star graph with root $i^* \in \mathcal{V}$ and $\mathcal{E} = \mathcal{N}(i^*)$, then

$$= \begin{cases} (\arg) \min_{\substack{\mu_{i^*, \cdot} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_{\cdot}) \\ \mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \langle f_{i^*, j}, \mu_{i^*, j} \rangle \\ \text{s.t. } \mu_{i^*} = \mu_{i^*, j} \circ \pi_{i^*}^{-1}, \quad \forall (i^*, j) \in \mathcal{N}(i^*). \end{cases} \quad (\text{star relaxation})$$

Proof.

- i. The first claim—the equivalence of the minimization of the sum of bivariate and its global relaxation—is a generic result, which can be found in [Theorem C.3](#).
- ii. The second claim—the equivalence of the minimization of the global relaxation and the tree relaxation of the sum of bivariate—can be proven by induction over trees \mathcal{G} : The fact is trivially true if \mathcal{G} has only one edge, as then the global measure and its 2-marginals are equivalent.

We prove now that if the second claim holds for all trees with $s \in \mathbb{N}$ edges, then it must also hold for all trees with $s + 1$ edges. Let \mathcal{G} now have $s + 1$ edges, i.e. $|\mathcal{E}| = s + 1$, and w.l.o.g. let (i^{j^*}, j^*) denote an edge, j^* be a leaf, then

$$\begin{aligned} (\arg) \min_{\mu \in \mathcal{M}_1^+(\Omega)} \langle F, \mu \rangle &= (\arg) \min_{\mu \in \mathcal{M}_1^+(\Omega)} \left\langle f_{k^*, i^{j^*}} + f_{i^{j^*}, j^*} + \sum_{\substack{(i, j) \in \mathcal{E} \\ (i, j) \neq (i^{j^*}, j^*) \\ (i, j) \neq (k^*, i^{j^*})}} f_{i, j}, \mu \right\rangle \\ &\quad (\text{as } G \text{ is a tree, there exists some } (k^*, i^{j^*}) \in \mathcal{E}) \\ &= (\arg) \min_{\mu \in \mathcal{M}_1^+(\Omega)} \left\langle f_{k^*, i^{j^*}} + m_{i^{j^*}, j^*} + \sum_{\substack{(i, j) \in \mathcal{E} \\ (i, j) \neq (i^{j^*}, j^*) \\ (i, j) \neq (k^*, i^{j^*})}} f_{i, j}, \mu \right\rangle \\ &\quad (\text{by } m_{i^{j^*}, j^*}(\cdot) := \min_{x_{j^*} \in \Omega_{j^*}} f_{i^{j^*}, j^*}(\cdot, x_{j^*}); j^* \text{ is a leaf}) \\ &= \begin{cases} (\arg) \min_{\substack{\mu_{\cdot, \cdot} \in \mathcal{M}^+(\Omega_{\cdot} \times \Omega_{\cdot}) \\ \mu_{\cdot} \in \mathcal{M}_1^+(\Omega_{\cdot})}} \langle f_{k^*, i^{j^*}} + m_{i^{j^*}, j^*}, \mu_{k^*, i^{j^*}} \rangle + \sum_{\substack{(i, j) \in \mathcal{E} \\ (i, j) \neq (i^{j^*}, j^*) \\ (i, j) \neq (k^*, i^{j^*})}} \langle f_{i, j}, \mu_{i, j} \rangle \\ \text{s.t. } \mu_i = \mu_{i, j} \circ \pi_i^{-1}, \quad \forall (i, j) \in \mathcal{E} \setminus \{(i^{j^*}, j^*)\} \\ \mu_j = \mu_{i, j} \circ \pi_j^{-1}, \quad \forall (i, j) \in \mathcal{E} \setminus \{(i^{j^*}, j^*)\} \end{cases} \\ &\quad (\text{induction hypothesis; } m_{i^{j^*}, j^*} : \Omega_{i^{j^*}} \rightarrow \mathbb{R}) \\ &= \begin{cases} (\arg) \min_{\substack{\mu_{\cdot, \cdot} \in \mathcal{M}^+(\Omega_{\cdot} \times \Omega_{\cdot}) \\ \mu_{\cdot} \in \mathcal{M}_1^+(\Omega_{\cdot})}} \langle m_{i^{j^*}, j^*}, \mu_{i^{j^*}} \rangle + \sum_{\substack{(i, j) \in \mathcal{E} \\ (i, j) \neq (i^{j^*}, j^*)}} \langle f_{i, j}, \mu_{i, j} \rangle \\ \text{s.t. } \mu_i = \mu_{i, j} \circ \pi_i^{-1}, \quad \forall (i, j) \in \mathcal{E} \setminus \{(i^{j^*}, j^*)\} \\ \mu_j = \mu_{i, j} \circ \pi_j^{-1}, \quad \forall (i, j) \in \mathcal{E} \setminus \{(i^{j^*}, j^*)\} \end{cases} \\ &\quad (m_{i^{j^*}, j^*} : \Omega_{i^{j^*}} \rightarrow \mathbb{R}) \\ &= \begin{cases} (\arg) \min_{\substack{\mu_{\cdot, \cdot} \in \mathcal{M}^+(\Omega_{\cdot} \times \Omega_{\cdot}) \\ \mu_{\cdot} \in \mathcal{M}_1^+(\Omega_{\cdot})}} \langle m_{i^{j^*}, j^*}, \mu_{i^{j^*}, j^*} \rangle + \sum_{\substack{(i, j) \in \mathcal{E} \\ (i, j) \neq (i^{j^*}, j^*)}} \langle f_{i, j}, \mu_{i, j} \rangle \\ \text{s.t. } \mu_i = \mu_{i, j} \circ \pi_i^{-1}, \quad \forall (i, j) \in \mathcal{E} \\ \mu_j = \mu_{i, j} \circ \pi_j^{-1}, \quad \forall (i, j) \in \mathcal{E} \end{cases} \\ &\quad (\text{introduce } \mu_{j^*} \in \mathcal{M}_1^+(\Omega_{j^*}), \mu_{i^{j^*}, j^*} \in \mathcal{M}^+(\Omega_{i^{j^*}} \times \Omega_{j^*}) : \mu_{i^{j^*}} = \mu_{i^{j^*}, j^*} \circ \pi_{i^{j^*}}^{-1}, \mu_{j^*} = \mu_{i^{j^*}, j^*} \circ \pi_{j^*}^{-1}) \end{aligned}$$

$$\begin{aligned}
&= \begin{cases} (\arg) \min_{\substack{\mu_{\cdot,\cdot} \in \mathcal{M}^+(\Omega \times \Omega) \\ \mu_{\cdot} \in \mathcal{M}_1^+(\Omega)}} \sum_{(i,j) \in \mathcal{E}} \langle f_{i,j}, \mu_{i,j} \rangle \\ \text{s.t. } \mu_i = \mu_{i,j} \circ \pi_i^{-1}, \quad \forall (i,j) \in \mathcal{E} \\ \mu_j = \mu_{i,j} \circ \pi_j^{-1}, \quad \forall (i,j) \in \mathcal{E}, \end{cases} \\
&\quad (\text{by } m_{ij^*,j^*}(\cdot) := \min_{x_{j^*} \in \Omega_{j^*}} f_{ij^*,j^*}(\cdot, x_{j^*}); j^* \text{ is a leaf})
\end{aligned}$$

which proves the desired result by induction.

- iii. The third claim—the equivalence of the tree relaxation and the star relaxation of the sum of bivariate—immediately follows from the fact that the unary measures of leafs are not part of the objective function and do not impose any non-trivial constraints. \square

Given we have marginally-coupled bivariate measures $\mu_{i,j}$ on $\Omega_i \times \Omega_j$, $(i,j) \in \mathcal{E}$, such as the candidates of the tree relaxation in [Theorem 3.3](#), we can ask the question of whether there exists a (global) measure on Ω that has these bivariate measures as marginals. [Theorem 3.4](#) proves existence of such a global measure for tree-indexed sets of bivariate measures and constructs one such global measure with maximum entropy.

Theorem 3.4 (Reconstruction From Marginals). *If $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is an oriented tree with root $i^* \in \mathcal{V}$ and given*

$$\begin{cases} \mu_{i,j} \in \mathcal{M}^+(\Omega_i \times \Omega_j), & \forall (i,j) \in \mathcal{E} \\ \mu_i \in \mathcal{M}_1^+(\Omega_i), & \forall i \in \mathcal{V} \\ \text{s.t. } \mu_i = \mu_{i,j} \circ \pi_i^{-1}, & \forall (i,j) \in \mathcal{E} \\ \mu_j = \mu_{i,j} \circ \pi_j^{-1}, & \forall (i,j) \in \mathcal{E}, \end{cases}$$

then the measure $\mu \in \mathcal{M}_1^+(\Omega)$ with

$$\mu(x) := \begin{cases} \mu_{i^*}(x_{i^*}) \prod_{(i,j) \in \mathcal{E}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} & \text{if } \mu_i(x_i) > 0 \quad \forall i \in \mathcal{V}, \quad \forall (x_1, \dots, x_n) \in \Omega \\ 0 & \text{else} \end{cases}$$

is a maximal-entropy measure, such that it has marginals $(\mu_{i,j}), (i,j) \in \mathcal{E}$, i.e.

$$\mu_{i,j} = \mu \circ \pi_{i,j}^{-1} \quad \text{for all } (i,j) \in \mathcal{E}.$$

Remark 2. The definition of μ in [Theorem 3.4](#) does not depend on the choice of the root of a tree. See [Corollary 3.1](#).

Proof. We assume w.l.o.g. that for all $x \in \Omega$ we have $\mu(x) > 0$, otherwise restrict Ω .

- i. First, we prove that μ is a probability measure. To this end, let (i^*, j^*) denote an edge and j^* be a leaf of \mathcal{G} . We have

$$\begin{aligned}
\mu(\Omega) &= \sum_{x \in \Omega} \mu(x) && \text{(additivity of measures)} \\
&= \sum_{\substack{x_i \in \Omega_i \\ i \neq j^*}} \sum_{x_{j^*} \in \Omega_{j^*}} \mu(x_1, \dots, x_n) && \text{(reorder summation)} \\
&= \sum_{\substack{x_i \in \Omega_i \\ i \neq j^*}} \sum_{x_{j^*} \in \Omega_{j^*}} \mu_{i^*}(x_{i^*}) \prod_{(i,j) \in \mathcal{E}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} && \text{(definition of } \mu)
\end{aligned}$$

$$\begin{aligned}
&= \sum_{\substack{x_i \in \Omega_i \\ i \neq j^*}} \mu_{i^*}(x_{i^*}) \left(\prod_{\substack{(i,j) \in \mathcal{E} \\ j \neq j^*}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} \right) \sum_{x_{j^*} \in \Omega_{j^*}} \frac{\mu_{ij^*,j^*}(x_{ij^*}, x_{j^*})}{\mu_{ij^*}(x_{ij^*})} \\
&\quad \text{(distributivity; } j^* \text{ is leaf)} \\
&= \sum_{\substack{x_i \in \Omega_i \\ i \neq j^*}} \mu_{i^*}(x_{i^*}) \left(\prod_{\substack{(i,j) \in \mathcal{E} \\ j \neq j^*}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} \right) \quad \text{(by } \mu_{ij^*} = \mu_{ij^*,j^*} \circ \pi_{ij^*}^{-1} \text{)} \\
&= \sum_{x_{i^*} \in \Omega_{i^*}} \mu_{i^*}(x_{i^*}) \quad \text{(repeated application of the last few steps for remaining leafs)} \\
&= 1. \quad \text{(by } \mu_{i^*} \in \mathcal{M}_1^*(\Omega_{i^*}) \text{)}
\end{aligned}$$

- ii. Similarly, we prove that $\mu \circ \pi_{i,j}^{-1} = \mu_{i,j}$ for all $(i,j) \in \mathcal{E}$. W.l.o.g., we assume that for an arbitrary edge $(i^*, j^*) \in \mathcal{E}$ the vertex j^* is a leaf and that there are no other leafs—otherwise can repeatedly marginalize w.r.t. leafs like in [Item i.](#). Denote with $(i^*, j^{i^*}) \in \mathcal{E}$ the edge that includes the root. We have for all $x_{ij^*} \in \Omega_{ij^*}, x_{j^*} \in \Omega_{j^*}$ that

$$\begin{aligned}
&\mu \circ \pi_{ij^*,j^*}^{-1}(x_{ij^*}, x_{j^*}) \\
&= \sum_{\substack{x_i \in \Omega_i \\ i \neq i^*, j^*}} \mu_{i^*}(x_{i^*}) \prod_{(i,j) \in \mathcal{E}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} \quad \text{(explicit summation; definition of } \mu \text{)} \\
&= \sum_{\substack{x_i \in \Omega_i \\ i \neq i^*, i^{j^*}, j^*}} \left(\sum_{\substack{x_{i^*} \in \Omega_{i^*} \\ \mu_{i^*}(x_{i^*}) > 0}} \frac{\mu_{i^*,j^{i^*}}(x_{i^*}, x_{j^{i^*}})}{\mu_{i^*}(x_{i^*})} \right) \prod_{\substack{(i,j) \in \mathcal{E} \\ i \neq i^*}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} \\
&\quad \text{(distributivity; } i^* \text{ is a root)} \\
&= \sum_{\substack{x_i \in \Omega_i \\ i \neq i^*, i^{j^*}, j^*}} \mu_{j^{i^*}}(x_{j^{i^*}}) \prod_{\substack{(i,j) \in \mathcal{E} \\ i \neq i^*}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} \quad \text{(cancellation; } \mu_{j^{i^*}} = \mu_{i^*,j^{i^*}} \circ \pi_{j^{i^*}}^{-1} \text{)} \\
&= \frac{\mu_{ij^*,j^*}(x_{ij^*}, x_{j^*})}{\mu_{ij^*}(x_{ij^*})}. \quad \text{(repeated application of the last few steps for remaining roots)}
\end{aligned}$$

- iii. Finally, we prove that μ is a unique maximal-entropy measure, such that it has marginals $(\mu_{i,j}), (i,j) \in \mathcal{E}$.

For an arbitrary measure $\lambda \in \mathcal{M}(\Omega)$ with $\mu_{ij} = \lambda \circ \pi_{ij}^{-1}$ for all $(i,j) \in \mathcal{E}$, we have

$$\begin{aligned}
H(\lambda) &= H(\lambda, \mu) - D_{\text{KL}}(\lambda, \mu) \quad \text{(cross-entropy formula)} \\
&= -\mathbb{E}_\lambda(\log \mu) - D_{\text{KL}}(\lambda, \mu) \quad \text{(definition of cross-entropy)} \\
&= -\mathbb{E}_\lambda \left(\log \mu_{i^*}(\cdot_{i^*}) \prod_{\substack{(i,j) \in \mathcal{E} \\ \mu_i(\cdot_i) > 0}} \frac{\mu_{i,j}(\cdot_i, \cdot_j)}{\mu_i(\cdot_i)} \right) - D_{\text{KL}}(\lambda, \mu) \quad \text{(definition of } \mu \text{)} \\
&= -\mathbb{E}_\lambda(\log \mu_{i^*}(\cdot_{i^*})) \\
&\quad - \left(\sum_{\substack{(i,j) \in \mathcal{E} \\ \mu_i(\cdot_i) > 0}} \mathbb{E}_\lambda(\log \mu_{i,j}(\cdot_i, \cdot_j)) - \mathbb{E}_\lambda(\log \mu_i(\cdot_i)) \right) - D_{\text{KL}}(\lambda, \mu) \\
&\quad \text{(properties of log and linearity of } \mathbb{E} \text{)} \\
&= -\mathbb{E}_{\lambda \circ \pi_{i^*}^{-1}}(\log \mu_{i^*})
\end{aligned}$$

$$\begin{aligned}
& - \left(\sum_{\substack{(i,j) \in \mathcal{E} \\ \mu_i(\cdot_i) > 0}} \mathbb{E}_{\lambda \circ \pi_{ij}^{-1}}(\log \mu_{i,j}) - \mathbb{E}_{\lambda \circ \pi_i^{-1}}(\log \mu_i(\cdot_i)) \right) - D_{\text{KL}}(\lambda, \mu) \\
& \quad \text{(expect. val. of functions constant in an argument)} \\
& = -\mathbb{E}_{\mu_{i^*}}(\log \mu_{i^*}) \\
& \quad - \left(\sum_{\substack{(i,j) \in \mathcal{E} \\ \mu_i(\cdot_i) > 0}} \mathbb{E}_{\mu_{ij}}(\log \mu_{i,j}(\cdot, \cdot)) - \mathbb{E}_{\mu_i}(\log \mu_i(\cdot_i)) \right) - D_{\text{KL}}(\lambda, \mu) \\
& \quad \text{(assumption } \mu_{ij} = \lambda \circ \pi_{ij}^{-1} \text{ for all } (i,j) \in \mathcal{E}) \\
& = -\mathbb{E}_{\mu}(\log \mu) - D_{\text{KL}}(\lambda, \mu) \quad \text{(by the reversed argument)} \\
& = H(\mu) - D_{\text{KL}}(\lambda, \mu) \quad \text{(definition of entropy)} \\
& \leq H(\mu). \quad \text{(property of the KL-divergence)}
\end{aligned}$$

□

It turns out that the limits to the existence of a global measure that has given marginally-coupled bivariate measures as marginals are related to the consistency of the tree relaxation of [Theorem 3.3](#) in solving the original minimization of sums of bivariate.

We see in [Example 3.5](#) that, in general, if we apply the tree relaxation to non-tree-structured sums of bivariate, we obtain a different problem. The problem we obtain differs both in its minimal value and in our ability to make sense of solutions by finding a global measure for a solution that is a set of marginally-coupled bivariate measures.

Example 3.5 (Extension Impossible). *Consider the functions $F : \{0, 1\}^3 \rightarrow \mathbb{R}$, $f_{1,2}, f_{1,3}, f_{2,3} : \{0, 1\}^2 \rightarrow \mathbb{R}$ with*

$$F(x_1, x_2, x_3) = f_{1,2}(x_1, x_2) + f_{1,3}(x_1, x_3) + f_{2,3}(x_2, x_3) \quad \forall x \in \{0, 1\}^3,$$

where

$$\begin{aligned}
f_{1,2}(x_1, x_2) &:= \begin{cases} 1 & \text{if } x_1 = x_2 = 0 \\ 0 & \text{if } x_1 = x_2 = 1 \\ \infty & \text{else,} \end{cases} \\
f_{1,3}(x_1, x_3) &:= \begin{cases} 1 & \text{if } x_1 = 1 - x_3 = 0 \\ 0 & \text{if } x_1 = x_3 \\ \infty & \text{else, and} \end{cases} \\
f_{2,3}(x_2, x_3) &:= \begin{cases} 1 & \text{if } x_2 = 1 - x_3 = 0 \\ 0 & \text{if } x_3 = 1 - x_2 = 0 \\ \infty & \text{else.} \end{cases}
\end{aligned}$$

Then, we have

$$F(x_1, x_2, x_3) = \begin{cases} 3 & \text{if } x_1 = x_2 = 1 - x_3 = 0 \\ \infty & \text{else} \end{cases}$$

and the tree relaxation of [Theorem 3.3](#) (where in this example $\infty \cdot 0 := 0$ and $\mathcal{G} := (\mathcal{V}, \mathcal{E}) := (\{1, 2, 3\}, \{(1, 2), (1, 3), (2, 3)\})$ is not a tree) yields

$$\begin{cases} (\arg) \min \sum_{\substack{\mu_{\cdot, \cdot} \in \mathcal{M}^+(\Omega \times \Omega) \\ \mu_{\cdot, \cdot} \in \mathcal{M}_1^+(\Omega)}} \langle f_{i,j}, \mu_{i,j} \rangle \\ \text{s.t. } \mu_i = \mu_{i,j} \circ \pi_i^{-1}, \quad \forall (i,j) \in \mathcal{E} \\ \mu_j = \mu_{i,j} \circ \pi_j^{-1}, \quad \forall (i,j) \in \mathcal{E} \end{cases}$$

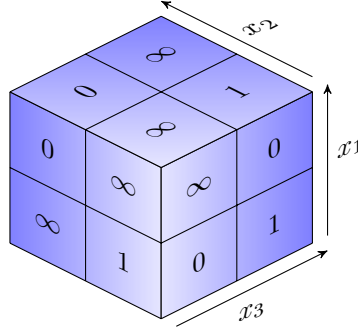


Figure 4: Geometric visualization of Example 3.5. Measures $\mu_{1,2}$, $\mu_{1,3}$, and $\mu_{2,3}$ can be understood to be defined on the respective face of the cube. The relaxed objective function is the sum of the scalar product of each measure with the respective face.

$$\begin{aligned}
 & \left\{ \begin{array}{l} (\arg) \min_{\mu_{\cdot,\cdot} \in \mathcal{M}^+(\Omega \times \Omega)} \langle f_{1,2}, \mu_{1,2} \rangle + \langle f_{1,3}, \mu_{1,3} \rangle + \langle f_{2,3}, \mu_{2,3} \rangle \\ s.t. \quad \mu_{1,2}(\Omega_1 \times \Omega_2) = 1 \\ \mu_{1,3}(\Omega_1 \times \Omega_3) = 1 \\ \mu_{2,3}(\Omega_2 \times \Omega_3) = 1 \\ \mu_{1,2} \circ \pi_1^{-1} = \mu_{1,3} \circ \pi_1^{-1} \\ \mu_{1,2} \circ \pi_2^{-1} = \mu_{2,3} \circ \pi_2^{-1} \\ \mu_{1,3} \circ \pi_3^{-1} = \mu_{2,3} \circ \pi_3^{-1} \end{array} \right. \\
 & \hspace{15em} (\text{definition of } \mathcal{E}; \text{ eliminating marginal variables}) \\
 & = \left\{ \begin{array}{l} (\arg) \min_{\mu_{\cdot,\cdot} \in \mathcal{M}^+(\Omega \times \Omega)} \mu_{1,2}(0,0) + \mu_{1,3}(0,1) + \mu_{2,3}(0,1) \\ s.t. \quad \mu_{1,2}(0,0) + \mu_{1,2}(1,0) + \mu_{1,2}(0,1) + \mu_{1,2}(1,1) = 1 \\ \mu_{1,3}(0,0) + \mu_{1,3}(1,0) + \mu_{1,3}(0,1) + \mu_{1,3}(1,1) = 1 \\ \mu_{2,3}(0,0) + \mu_{2,3}(1,0) + \mu_{2,3}(0,1) + \mu_{2,3}(1,1) = 1 \\ \mu_{1,2}(0,0) + \mu_{1,2}(0,1) = \mu_{1,3}(0,0) + \mu_{1,3}(0,1) \\ \mu_{1,2}(1,0) + \mu_{1,2}(1,1) = \mu_{1,3}(1,0) + \mu_{1,3}(1,1) \\ \mu_{1,2}(0,0) + \mu_{1,2}(1,0) = \mu_{2,3}(0,0) + \mu_{2,3}(0,1) \\ \mu_{1,2}(0,1) + \mu_{1,2}(1,1) = \mu_{2,3}(1,0) + \mu_{2,3}(1,1) \\ \mu_{1,3}(0,0) + \mu_{1,3}(1,0) = \mu_{2,3}(0,0) + \mu_{2,3}(1,0) \\ \mu_{1,3}(0,1) + \mu_{1,3}(1,1) = \mu_{2,3}(0,1) + \mu_{2,3}(1,1) \\ \hline \mu_{1,2}(0,1) = \mu_{1,2}(1,0) = \mu_{1,3}(1,0) = \mu_{2,3}(0,0) = \mu_{2,3}(1,1) = 0 \end{array} \right. \\
 & \hspace{15em} (\text{definitions of } f_{\cdot,\cdot} \text{ \& } \Omega_{\cdot})
 \end{aligned}$$

$$\begin{aligned}
&= \left\{ \begin{array}{l} (\arg) \min_{\mu_{\cdot,\cdot} \in \mathcal{M}^+(\Omega \times \Omega)} \mu_{1,2}(0,0) + \mu_{1,3}(0,1) + \mu_{2,3}(0,1) \\ s.t. \quad \mu_{1,2}(0,0) + \mu_{1,2}(1,1) = 1 \\ \mu_{1,3}(0,0) + \mu_{1,3}(0,1) + \mu_{1,3}(1,1) = 1 \\ \mu_{2,3}(1,0) + \mu_{2,3}(0,1) = 1 \\ \mu_{1,2}(0,0) = \mu_{1,3}(0,0) + \mu_{1,3}(0,1) \\ \mu_{1,2}(1,1) = \mu_{1,3}(1,1) \\ \mu_{1,2}(0,0) = \mu_{2,3}(0,1) \\ \mu_{1,2}(1,1) = \mu_{2,3}(1,0) \\ \mu_{1,3}(0,0) = \mu_{2,3}(1,0) \\ \mu_{1,3}(0,1) + \mu_{1,3}(1,1) = \mu_{2,3}(0,1) \end{array} \right. \\
&\quad \mu_{1,2}(0,1) = \mu_{1,2}(1,0) = \mu_{1,3}(1,0) = \mu_{2,3}(0,0) = \mu_{2,3}(1,1) = 0 \\
&\quad \text{(plugging the last equation into all others)} \\
&= \left\{ \begin{array}{l} (\arg) \min_{\mu_{\cdot,\cdot} \in \mathcal{M}^+(\Omega \times \Omega)} 2\mu_{1,2}(0,0) + \mu_{1,3}(0,1) \\ s.t. \quad 2\mu_{1,2}(1,1) + \mu_{1,3}(0,1) = 1 \\ \mu_{1,2}(0,0) = \mu_{1,2}(1,1) + \mu_{1,3}(0,1) \\ \mu_{1,3}(0,1) + \mu_{1,2}(1,1) = \mu_{1,2}(0,0) \end{array} \right. \\
&\quad \mu_{1,2}(0,1) = \mu_{1,2}(1,0) = \mu_{1,3}(1,0) = \mu_{2,3}(0,0) = \mu_{2,3}(1,1) = 0 \\
&\quad \mu_{1,3}(1,1) = \mu_{1,2}(1,1) \\
&\quad \mu_{2,3}(0,1) = \mu_{1,2}(0,0) \\
&\quad \mu_{1,3}(0,0) = \mu_{2,3}(1,0) = \mu_{1,2}(1,1) \\
&\quad \mu_{1,2}(1,1) = 1 - \mu_{1,2}(0,0) \\
&\quad \mu_{2,3}(1,0) = 1 - \mu_{1,2}(0,0) \\
&\quad \text{(replacing variables constrained by trivial equations)} \\
&= \left\{ \begin{array}{l} (\arg) \min_{\mu_{\cdot,\cdot} \in \mathcal{M}^+(\Omega \times \Omega)} 1 \\ s.t. \quad \mu_{1,2}(0,0) = 1 - \mu_{1,2}(1,1) \\ 1 - \mu_{1,2}(1,1) = \mu_{1,2}(0,0) \end{array} \right. \\
&\quad \mu_{1,2}(0,1) = \mu_{1,2}(1,0) = \mu_{1,3}(1,0) = \mu_{2,3}(0,0) = \mu_{2,3}(1,1) = 0 \\
&\quad \mu_{1,3}(1,1) = \mu_{1,2}(1,1) \\
&\quad \mu_{2,3}(0,1) = \mu_{1,2}(0,0) \\
&\quad \mu_{1,3}(0,0) = \mu_{2,3}(1,0) = \mu_{1,2}(1,1) \\
&\quad \mu_{1,2}(1,1) = 1 - \mu_{1,2}(0,0) \\
&\quad \mu_{2,3}(1,0) = 1 - \mu_{1,2}(0,0) \\
&\quad \mu_{1,3}(0,1) = 1 - 2\mu_{1,2}(1,1) . \\
&\quad \text{(replacing variables constrained by trivial equations)}
\end{aligned}$$

Clearly, we have a minimal value of 1 at a (non-unique) minimum of

$$\mu_{1,2}(a,b) = \mu_{1,3}(a,b) = 1/2 - \mu_{2,3}(a,b) = \begin{cases} 1/2 & \text{if } a = b \\ 0 & \text{else} \end{cases} \quad \forall a, b \in \{0,1\}.$$

Further, there exists no measure μ on $\{0,1\}^3$ such that it has marginals $\mu_{1,2}, \mu_{1,3}, \mu_{2,3}$,

i.e. there exists no μ such that

$$\mu_{1,2} = \mu \circ \pi_{1,2}^{-1}, \quad \mu_{1,3} = \mu \circ \pi_{1,3}^{-1}, \quad \text{and} \quad \mu_{2,3} = \mu \circ \pi_{2,3}^{-1}.$$

This is, as $\mu_{1,2}(1,0) = \mu_{1,2}(0,1) = \mu_{1,3}(1,0) = \mu_{1,3}(0,1) = 0$ would imply by marginalization that the support of μ is contained in $\{(0,0,0), (1,1,1)\}$. However, since $\mu_{2,3}(0,0) = \mu_{2,3}(1,1) = 0$, by marginalization, μ must be zero on $\{(\cdot,0,0), (\cdot,1,1)\}$. Therefore, μ cannot have the specified marginals.

This example shows that neither the non-global relaxation of [Theorem 3.3](#) nor the reconstruction of [Theorem 3.4](#) can be extended beyond trees in general.

The following two postulates are inspired by and very similar to the results in [\[OP24\]](#). The goal is to derive the Lagrangian dual of the star relaxation in [Postulate 3.1](#), which restricts to star graph-indexed sums of bivariates. The dual of the tree relaxation would follow similarly and will be described briefly in [Theorem 4.1](#) of [Section 4](#). Later in our derivation, even the specific results for star graph-indexed sums of bivariates will prove useful, as we will be able to use them for a type of coordinate ascent minimization algorithm, similar and motivated to the results in [\[OP24\]](#).

Further, [Postulate 3.1](#) claims that the dual of the star relaxation is strong and involves fewer variables than its primal—the star relaxation. This dual maximization can be understood as the objective to find lower bounds for the functions $(\min_{x_j \in \Omega_j} f_{i^*,j}(\cdot, x_j)) \in \mathbb{R}^{\Omega_{i^*}}$ for all $(i^*, j) \in \mathcal{N}(i^*)$, called min-marginals, such that the sum of the min-marginals has the largest possible minimum.

Postulate 3.1 (Star Lagrangian Duality). *If $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a star graph with root $i^* \in \mathcal{V}$ and $\mathcal{E} = \mathcal{N}(i^*)$, the Lagrangian dual of the star relaxation of [Theorem 3.3](#) is strong and*

$$\begin{aligned} & \begin{cases} \min_{\substack{\mu_{i^*,\cdot} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_{\cdot}) \\ \mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})}} \sum_{(i^*,j) \in \mathcal{N}(i^*)} \langle f_{i^*,j}, \mu_{i^*,j} \rangle \\ \text{s.t. } \mu_{i^*} = \mu_{i^*,j} \circ \pi_{i^*}^{-1}, \quad \forall (i^*,j) \in \mathcal{N}(i^*) \end{cases} & \text{(star relaxation)} \\ & = \begin{cases} \max_{\rho_{i^*,\cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*,j) \in \mathcal{N}(i^*)} \rho_{i^*,j}(x_{i^*}) \right) \right) \\ \text{s.t. } \rho_{i^*,j} \leq \min_{x_j \in \Omega_j} f_{i^*,j}(\cdot, x_j), \quad \forall (i^*,j) \in \mathcal{N}(i^*). \end{cases} & \text{(dual star relaxation)} \end{aligned}$$

We conclude the analysis of the star relaxation by deriving a closed-form solution to the dual star relaxation in [Postulate 3.2](#) using [Lemma 3.1](#). The solution is not unique, but is characterized by linear inequalities. Intuitively, optimal solutions exceed the minimum of the sum of the min-marginals, while being individually bounded by them.

Lemma 3.1. *In the setting of [Definition 1.1](#) and [Postulate 3.1](#), we have for all $(i^*, j) \in \mathcal{N}(i^*)$ that*

$$\begin{aligned} \min_{\mu_{i^*,j} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_j)} \langle f_{i^*,j} - \rho_{i^*,j}, \mu_{i^*,j} \rangle &= \min_{\nu_{i^*,j} \in \mathcal{M}^+(\Omega_{i^*})} \langle m_{i^*,j} - \rho_{i^*,j}, \nu_{i^*,j} \rangle \\ &= \delta_{\rho_{i^*,j} \leq m_{i^*,j}} := \begin{cases} 0 & \text{if } \rho_{i^*,j} \leq m_{i^*,j} \\ -\infty & \text{else,} \end{cases} \end{aligned}$$

where $m_{i^*,j}(\cdot) := \min_{x_j \in \Omega_j} f_{i^*,j}(\cdot, x_j)$ and $\rho_{i^*,j} \in \mathbb{R}^{\Omega_{i^*}}$.

Proof. Clear. □

Postulate 3.2 (Necessary and Sufficient Conditions for Star Dual Variables). *If $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a star graph with root $i^* \in \mathcal{V}$ and $\mathcal{E} = \mathcal{N}(i^*)$, the variables of the Lagrangian dual of the star relaxation of [Postulate 3.1](#) are optimal, i.e.*

$$\begin{aligned} \rho_{i^*, \cdot}^* &\in \begin{cases} \arg \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right) (x_{i^*}) \right) \\ \text{s.t.} \quad \rho_{i^*, j} \leq m_{i^*, j}, \quad \forall (i^*, j) \in \mathcal{N}(i^*) \end{cases} & \text{(dual star relaxation)} \\ \iff &\begin{cases} \min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} m_{i^*, j} \right) (x_{i^*}) \leq \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \\ \text{s.t.} \quad \rho_{i^*, j} \leq m_{i^*, j}, \quad \forall (i^*, j) \in \mathcal{N}(i^*), \end{cases} & \text{(NSCSDV)} \end{aligned}$$

where $m_{i^*, j}(\cdot) := \min_{x_j \in \Omega_j} f_{i^*, j}(\cdot, x_j)$ for all $(i^*, j) \in \mathcal{N}(i^*)$.

3.3. Entropy-Regularized Relaxation Principle

As solutions to the dual star relaxation from [Postulate 3.2](#) are not unique, and with applications of the results to non-tree-structured sums of bivariate marginals in mind, we aim to derive an *entropy-regularized star relaxation*. The entropy-regularization will, in fact, turn out to help us in deriving a unique solution, and the new objective that we will construct will uniformly approximate our known star relaxation. The heuristic hope is that an entropy-regularization would also help in “distributing” measure by preferring higher entropies in coupled bivariate marginal measures, such that the possibility of suboptimal convergence may be partly prevented.

To this end, we formulate a “parameter-free” version of the maximal-entropy measure that has given coupled bivariate marginal measures, known from [Theorem 3.4](#), in [Corollary 3.1](#). We also find a simple expression for its entropy that is a linear combination of the entropies of its marginals.

Corollary 3.1 (Representation of the Maximal-Entropy Reconstruction). *If the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a tree and given*

$$\begin{cases} \mu_{i,j} \in \mathcal{M}^+(\Omega_i \times \Omega_j), & \forall (i,j) \in \mathcal{E} \\ \mu_i \in \mathcal{M}_1^+(\Omega_i), & \forall i \in \mathcal{V} \\ \text{s.t.} \quad \mu_i = \mu_{i,j} \circ \pi_i^{-1}, & \forall (i,j) \in \mathcal{E} \\ \mu_j = \mu_{i,j} \circ \pi_j^{-1}, & \forall (i,j) \in \mathcal{E}, \end{cases}$$

then the maximal-entropy measure $\mu \in \mathcal{M}_1^+(\Omega)$ with marginals $(\mu_{i,j}), (i,j) \in \mathcal{E}$, i.e.

$$\mu_{i,j} = \mu \circ \pi_{i,j}^{-1} \quad \forall (i,j) \in \mathcal{E},$$

known from [Theorem 3.4](#), has a representation

$$\mu(x) = \begin{cases} \left(\prod_{(i,j) \in \mathcal{E}} \mu_{i,j}(x_i, x_j) \right) / \left(\prod_{i \in \mathcal{V}} \mu_i(x_i)^{|\tilde{\mathcal{N}}(i)|-1} \right) & \text{if } \mu_i(x_i) > 0 \quad \forall i \in \mathcal{V}, \quad \forall x \in \Omega. \\ 0 & \text{else} \end{cases}$$

Further, we have for its entropy

$$H(\mu) = \sum_{(i,j) \in \mathcal{E}} H(\mu_{i,j}) - \sum_{i \in \mathcal{V}} (|\tilde{\mathcal{N}}(i)| - 1) \cdot H(\mu_i).$$

Proof. W.l.o.g., we assume \mathcal{G} to be oriented. Otherwise, we orient it. Note, that we now have

$$|\tilde{\mathcal{N}}(i)| - 1 = |\mathcal{N}(i)|, \quad \forall i \in \mathcal{V} \setminus \{i^*\},$$

where the root of the oriented tree is denoted by $i^* \in \mathcal{V}$.

- i. We have to prove that μ as defined in [Theorem 3.4](#) equals the representation given. Clearly, this is the case if there exists $x \in \Omega$ and $i \in \mathcal{V}$ with $\mu_i(x_i) = 0$ by the “else”-condition of the definitions. Therefore, we assume that for all $x \in \Omega$, we have $\mu_i(x_i) > 0$ for all $i \in \mathcal{V}$.

We prove the equality by induction over trees. The fact is trivially true for a tree with the root as the only parent vertex. We assume, therefore, that the equality holds for trees of a fixed depth. Denote by $(i^{j^*}, j^*) \in \mathcal{E}$ an arbitrary leaf edge of \mathcal{G} . Then,

$$\begin{aligned} \mu(x) &= \mu_{i^*}(x_{i^*}) \prod_{(i,j) \in \mathcal{E}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} && \text{(definition in Theorem 3.4)} \\ &= \left(\mu_{i^*}(x_{i^*}) \prod_{(i,j) \in \mathcal{E} \setminus \{(i^{j^*}, j^*)\}} \frac{\mu_{i,j}(x_i, x_j)}{\mu_i(x_i)} \right) \frac{\mu_{i^{j^*}, j^*}(x_{i^{j^*}}, x_{j^*})}{\mu_{i^{j^*}}(x_{i^{j^*}})} \\ & && \text{(splitting the product)} \\ &= \left(\left(\prod_{(i,j) \in \mathcal{E} \setminus \{(i^{j^*}, j^*)\}} \mu_{i,j}(x_i, x_j) \right) / \left(\prod_{i \in (\mathcal{V} \setminus \{j^*\})} \mu_i(x_i)^{|\mathcal{N}^{j^*}(i)|} \right) \right) \frac{\mu_{i^{j^*}, j^*}(x_{i^{j^*}}, x_{j^*})}{\mu_{i^{j^*}}(x_{i^{j^*}})} \\ & && \text{(induction hypothesis; } \mathcal{N}^{j^*}(i) := \{(i, j) \in \mathcal{E} \setminus \{(i^{j^*}, j^*)\}\}) \\ &= \left(\prod_{(i,j) \in \mathcal{E}} \mu_{i,j}(x_i, x_j) \right) / \left(\prod_{i \in \mathcal{V}} \mu_i(x_i)^{|\mathcal{N}^{j^*}(i)|} \right). \\ & && \text{(by } |\mathcal{N}^{j^*}(i)| = |\mathcal{N}(i)|, \text{ if } i \neq i^{j^*}; \quad |\mathcal{N}^{j^*}(i^{j^*})| + 1 = |\mathcal{N}(i^{j^*})|; \quad |\mathcal{N}(j^*)| = 0) \end{aligned}$$

- ii. W.l.o.g. assume $\mu(x) > 0$ for all $x \in \Omega$. Otherwise, restrict Ω . We have

$$\begin{aligned} H(\mu) &= -\mathbb{E}_\mu(\log \mu(\cdot)) && \text{(definition of the entropy } H) \\ &= -\mathbb{E}_\mu \left(\log \left(\prod_{(i,j) \in \mathcal{E}} \mu_{i,j}(\cdot_i, \cdot_j) \right) / \left(\prod_{i \in \mathcal{V}} \mu_i(\cdot_i)^{|\tilde{\mathcal{N}}(i)|-1} \right) \right) && \text{(representation of } \mu) \\ &= - \sum_{(i,j) \in \mathcal{E}} \mathbb{E}_\mu \left(\log \mu_{i,j}(\cdot_i, \cdot_j) \right) + \sum_{i \in \mathcal{V}} (|\tilde{\mathcal{N}}(i)| - 1) \cdot \mathbb{E}_\mu \left(\log \mu_i(\cdot_i) \right) \\ & && \text{(properties of the logarithm; linearity of } \mathbb{E}) \\ &= - \sum_{(i,j) \in \mathcal{E}} \mathbb{E}_{\mu \circ \pi_{ij}^{-1}} \left(\log \mu_{i,j} \right) + \sum_{i \in \mathcal{V}} (|\tilde{\mathcal{N}}(i)| - 1) \cdot \mathbb{E}_{\mu \circ \pi_i^{-1}} \left(\log \mu_i \right) \\ & && \text{(expect. val. of functions constant in an argument)} \\ &= - \sum_{(i,j) \in \mathcal{E}} \mathbb{E}_{\mu_{ij}} \left(\log \mu_{i,j} \right) + \sum_{i \in \mathcal{V}} (|\tilde{\mathcal{N}}(i)| - 1) \cdot \mathbb{E}_{\mu_i} \left(\log \mu_i \right) \\ & && \text{(assumption } \mu_{ij} = \lambda \circ \pi_{ij}^{-1} \text{ for all } (i, j) \in \mathcal{E}) \\ &= \sum_{(i,j) \in \mathcal{E}} H(\mu_{ij}) - \sum_{i \in \mathcal{V}} (|\tilde{\mathcal{N}}(i)| - 1) \cdot H(\mu_i). && \text{(definition of the entropy } H) \end{aligned}$$

□

Next, we prove the concavity of the function that, given a family of bivariate measures, assigns the entropy of a global measure that has the family as marginals. This result will be essential to prove the strong Lagrangian duality in [Theorem 3.5](#).

Lemma 3.2. *If $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a star graph with root $i^* \in \mathcal{V}$ and $\mathcal{E} = \mathcal{N}(i^*)$, the function*

$$\mu_{i^*,j} \in \mathcal{M}_1^+(\Omega_{i^*} \times \Omega_j), (i^*, j) \in \mathcal{E} \mapsto \left(\sum_{(i^*,j) \in \mathcal{E}} H(\mu_{i^*,j}) \right) - (|\tilde{\mathcal{N}}(i^*)| - 1) \cdot H(\mu_{i^*,j} \circ \pi_{i^*}^{-1})$$

is concave.

Proof. Let $\mu_{i^*,j}, \lambda_{i^*,j} \in \mathcal{M}_1^+(\Omega_{i^*} \times \Omega_j)$, $t \in [0, 1]$, and define $z_{i^*,j} := t\mu_{i^*,j} + (1-t)\lambda_{i^*,j}$. Due to concavity of the entropy and $|\mathcal{E}| = |\tilde{\mathcal{N}}(i^*)|$, it is sufficient to show the concavity of

$$\begin{aligned} & \sum_{(i^*,j) \in \mathcal{E}} H(z_{i^*,j}) - H(z_{i^*,j} \circ \pi_{i^*}^{-1}) \\ &= \sum_{(i^*,j) \in \mathcal{E}} \log(|\Omega_j|) - D_{\text{KL}}(z_{i^*,j}, z_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j|) \quad (\text{see below argument}) \\ &= \sum_{(i^*,j) \in \mathcal{E}} \log(|\Omega_j|) \\ & \quad - D_{\text{KL}}\left(t\mu_{i^*,j} + (1-t)\lambda_{i^*,j}, t\mu_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j| + (1-t)\lambda_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j|\right) \\ & \quad (\text{definition of } z_{i^*,j} \text{ \& linearity of marginalization}) \\ &\geq \sum_{(i^*,j) \in \mathcal{E}} t \left(\log(|\Omega_j|) - D_{\text{KL}}\left(\mu_{i^*,j}, \mu_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j|\right) \right) \\ & \quad + (1-t) \left(\log(|\Omega_j|) - D_{\text{KL}}\left(\lambda_{i^*,j}, \lambda_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j|\right) \right) \\ & \quad (\text{convexity of KL-divergence}) \\ &= t \left(\sum_{(i^*,j) \in \mathcal{E}} H(\mu_{i^*,j}) - H(\mu_{i^*,j} \circ \pi_{i^*}^{-1}) \right) + (1-t) \left(\sum_{(i^*,j) \in \mathcal{E}} H(\lambda_{i^*,j}) - H(\lambda_{i^*,j} \circ \pi_{i^*}^{-1}) \right). \\ & \quad (\text{see below argument; linearity}) \end{aligned}$$

It follows a similar argument as in [\[GJ07, Appendix A\]](#). We have a representation of what is known as the conditional entropy of $\mu_{i^*,j}$, as

$$\begin{aligned} & H(\mu_{i^*,j}) - H(\mu_{i^*,j} \circ \pi_{i^*}^{-1}) \\ &= -\mathbb{E}_{\mu_{i^*,j}}(\log \mu_{i^*,j}) + \mathbb{E}_{\mu_{i^*,j} \circ \pi_{i^*}^{-1}}(\log \mu_{i^*,j} \circ \pi_{i^*}^{-1}) \quad (\text{definition entropy}) \\ &= -\mathbb{E}_{\mu_{i^*,j}}\left(\log \mu_{i^*,j} - \log \mu_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})\right) \quad (\text{constant integration}) \\ &= -\mathbb{E}_{\mu_{i^*,j}}\left(\log \frac{\mu_{i^*,j}}{\mu_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})}\right) \quad (\text{properties of log}) \\ &= \log(|\Omega_j|) - \log(|\Omega_j|) - \mathbb{E}_{\mu_{i^*,j}}\left(\log \frac{\mu_{i^*,j}}{\mu_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j|}\right) \quad (\text{constant addition}) \\ &= \log(|\Omega_j|) - \mathbb{E}_{\mu_{i^*,j}}\left(\log \frac{\mu_{i^*,j}(\cdot, \cdot)}{\mu_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j|}\right) \quad (\text{linearity \& properties of log}) \\ &= \log(|\Omega_j|) - D_{\text{KL}}(\mu_{i^*,j}, \mu_{i^*,j} \circ \pi_{i^*}^{-1}(\cdot_{i^*})/|\Omega_j|). \end{aligned}$$

(definition of KL-divergence; marginalized distribution is constant in marginalized coordinate)

□

By subtracting a scaled version of the entropy function of [Lemma 3.2](#) from the star relaxation of [Theorem 3.3](#), one obtains a new regularized objective, termed *entropy-regularized star relaxation*, which, due to the boundedness of the entropy, uniformly approximates the star relaxation. [Theorem 3.5](#) shows that the dual of this new objective now has a unique solution and derives its representation based on the ε -LogSumExp.

Similarly to the results in [Section 3.2](#), the theorem will also aid in deriving optimization algorithms that are applicable to generic sums of bivariates.

Theorem 3.5 (Entropy Star Lagrangian Duality). *If $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a star graph with root $i^* \in \mathcal{V}$ and $\mathcal{E} = \mathcal{N}(i^*)$, the Lagrangian dual of the following entropy-regularized star relaxation of [Theorem 3.3](#) is strong, has a unique solution, and*

$$\begin{aligned}
& \left\{ \begin{array}{l} \min_{\substack{\mu_{i^*, \cdot} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_{\cdot}) \\ \mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \langle f_{i^*, j}, \mu_{i^*, j} \rangle - \varepsilon H(\mu_{i^*, j}) \right) + \varepsilon(|\mathcal{N}(i^*)| - 1) \cdot H(\mu_{i^*}) \\ \text{s.t.} \quad \mu_{i^*} = \mu_{i^*, j} \circ \pi_{i^*}^{-1}, \quad \forall (i^*, j) \in \mathcal{N}(i^*) \end{array} \right. \quad (\text{entropy-regularized star relaxation}) \\
& = \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\left(\min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}(x_{i^*}) \right) \right. \right. \\
& \quad \left. \left. - \varepsilon \sum_{(i^*, j) \in \mathcal{N}(i^*)} \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{x_j \in \Omega_j} \exp((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j))/\varepsilon - 1) \right) \right) \\
& \quad (\text{entropy dual star relaxation}) \\
& = \left\{ \begin{array}{l} \rho - |\mathcal{N}(i^*)| \\ \text{where} \\ \rho_{i^*, j}(x_{i^*}) = \text{lse}_{\varepsilon}(f_{i^*, j}(x_{i^*}, \cdot)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_{\varepsilon}(f_{i^*, k}(x_{i^*}, \cdot)) \right) \\ \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*) \\ \rho = \varepsilon |\mathcal{N}(i^*)| + \varepsilon |\mathcal{N}(i^*)| \log(|\mathcal{N}(i^*)|/\varepsilon) \\ - |\mathcal{N}(i^*)| \text{lse}_{\varepsilon}^{x_{i^*}, j} \left(\text{lse}_{\varepsilon}^{x_j}(f_{i^*, j}(x_{i^*}, x_j)) + \text{lse}_{\varepsilon}^{x_j}(-f_{i^*, j}(x_{i^*}, x_j)) \right. \\ \quad \left. - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_{\varepsilon}^{x_k}(f_{i^*, k}(x_{i^*}, x_k)) \right) \end{array} \right.
\end{aligned}$$

Proof. The equivalence of the value of Lagrangian primal and dual, i.e. the strong duality, as well as the uniqueness of the solution can be derived directly, as

$$\begin{aligned}
& \left\{ \begin{array}{l} \min_{\substack{\mu_{i^*, \cdot} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_{\cdot}) \\ \mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \langle f_{i^*, j}, \mu_{i^*, j} \rangle - \varepsilon H(\mu_{i^*, j}) \right) + \varepsilon (|\mathcal{N}(i^*)| - 1) \cdot H(\mu_{i^*}) \\ \text{s.t. } \mu_{i^*} = \mu_{i^*, j} \circ \pi_{i^*}^{-1}, \quad \forall (i^*, j) \in \mathcal{N}(i^*). \end{array} \right. \\
&= \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \min_{\substack{\mu_{i^*, \cdot} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_{\cdot}) \\ \mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \langle f_{i^*, j}, \mu_{i^*, j} \rangle - \varepsilon H(\mu_{i^*, j}) \right) \\
&\quad + \varepsilon (|\mathcal{N}(i^*)| - 1) \cdot H(\mu_{i^*}) \\
&\quad + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \langle \rho_{i^*, j}, \mu_{i^*} - \mu_{i^*, j} \circ \pi_{i^*}^{-1} \rangle \\
&\quad \text{(strong Lagrangian duality by Theorem C.2; Lemma 3.2)} \\
&= \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \min_{\substack{\mu_{i^*, \cdot} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_{\cdot}) \\ \mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})}} \left\langle -\varepsilon (|\mathcal{N}(i^*)| - 1) \log \mu_{i^*} + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}, \mu_{i^*} \right\rangle \\
&\quad + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \left\langle \varepsilon \log \mu_{i^*, j} + f_{i^*, j} - \rho_{i^*, j}, \mu_{i^*, j} \right\rangle \\
&\quad \text{(definition entropy; linearity; constant integration)} \\
&= \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\min_{\mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})} \left\langle -\varepsilon (|\mathcal{N}(i^*)| - 1) \log \mu_{i^*} + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}, \mu_{i^*} \right\rangle \right) \\
&\quad + \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \min_{\mu_{i^*, j} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_j)} \left\langle \varepsilon \log \mu_{i^*, j} + f_{i^*, j} - \rho_{i^*, j}, \mu_{i^*, j} \right\rangle \right) \\
&\quad \text{(minimize separable problems)} \\
&= \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\left(\min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right)(x_{i^*}) \right) \right. \\
&\quad \left. - \varepsilon \sum_{(i^*, j) \in \mathcal{N}(i^*)} \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{x_j \in \Omega_j} \exp \left((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j)) / \varepsilon - 1 \right) \right) \\
&\quad \text{(by Lemma A.1 & Lemma A.2)} \\
&= \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\left(\min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right)(x_{i^*}) \right) \right. \\
&\quad \left. - \varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp (\rho_{i^*, j}(x_{i^*}) / \varepsilon) \sum_{x_j \in \Omega_j} \exp (-f_{i^*, j}(x_{i^*}, x_j) / \varepsilon - 1) \right) \\
&\quad \text{(distributivity & order of summation)}
\end{aligned}$$

$$= \begin{cases} \rho - |\mathcal{N}(i^*)| \\ \text{where} \\ \rho_{i^*,j}(x_{i^*}) = \text{lse}_\varepsilon(f_{i^*,j}(x_{i^*}, \cdot)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*,k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon(f_{i^*,k}(x_{i^*}, \cdot)) \right), \\ \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*) \\ \rho = \varepsilon |\mathcal{N}(i^*)| + \varepsilon |\mathcal{N}(i^*)| \log(|\mathcal{N}(i^*)|/\varepsilon) \\ - |\mathcal{N}(i^*)| \text{lse}_\varepsilon^{x_{i^*},j} \left(\text{lse}_\varepsilon^{x_j}(f_{i^*,j}(x_{i^*}, x_j)) + \text{lse}_\varepsilon^{x_j}(-f_{i^*,j}(x_{i^*}, x_j)) \right. \\ \left. - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*,k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon^{x_k}(f_{i^*,k}(x_{i^*}, x_k)) \right). \end{cases}$$

(by Lemma A.3)

□

3.4. Recovering Primal Solutions from Dual Solutions

Finally, we describe how to obtain a minimum of the sum of bivariates from a dual solution, e.g., in the problem formulations of [Postulates 3.1](#) and [3.2](#), [Theorems 3.5](#) and [4.1](#), and [Remark 4](#). Since we apply the method beyond star graphs, as opposed to much of the theory of this section, we prove that the result of [Theorem 3.6](#) holds for tree-indexed sums of bivariates. However, the result will be heuristically applicable beyond tree-indexed sums of bivariates.

As we have not introduced the dual of the tree relaxation explicitly for this case yet—mainly because closed-form solutions do not exist in this generality—we forward-reference to the description in [Theorem 4.1](#) to avoid redundancy.

Intuitively, the method in [Theorem 3.6](#) determines a solution by sequentially minimizing a sum which fixes determined arguments in bivariates, ignores bivariates that are constant in the active argument, and replaces all bivariates with two undertermined arguments by their min-marginals.

Theorem 3.6 (Recovering Primal Solutions). *In the setting of [Theorem 4.1](#), let*

- $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be an oriented tree,
- $\rho_{i,j} \in \mathbb{R}^{\Omega_i}, \rho_{j,i} \in \mathbb{R}^{\Omega_j}; (i,j) \in \mathcal{E}$ be a solution to the dual linear program,
- $i_j \in \mathcal{V}$ denote the parent node for all non-root $j \in \mathcal{V}$, and

$$m_{i,j}(\cdot) := \min_{x_j \in \Omega_j} f_{i,j}(\cdot, x_j) - \rho_{j,i}(x_j) \quad \forall (i,j) \in \mathcal{E}.$$

Then, we know that

$$x_j^* \in \arg \min_{x_j \in \Omega_{i_t}} \begin{cases} \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i_j\}} m_{j,k}(x_j) & \text{if } j \text{ is the root} \\ f_{i_j,j}(x_{i_j}^*, x_j) + \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i_j\}} m_{j,k}(x_j) & \text{else} \end{cases}$$

minimizes the sum of bivariates $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j} : \Omega \rightarrow \mathbb{R}$.

Proof. Define $\rho := \sum_{i \in \mathcal{V}} \rho_i \in \mathbb{R}$. By strong duality (similar proof as [Postulate 3.1](#)), we know $\min F - \rho = 0$ and that

$$F - \rho = \sum_{(i,j) \in \mathcal{E}} f_{i,j} - \rho_{i,j} - \rho_{j,i}.$$

By the constraint of the *dual linear program* of [Theorem 4.1](#), we know that each summand $f_{i,j} - \rho_{i,j} - \rho_{j,i}; (i,j) \in \mathcal{E}$ is non-negative, and, therefore, that

$$x^* \in \arg \min_{x \in \Omega} F(x) - \rho \iff \left(f_{i,j}(x_i^*, x_j^*) - \rho_{i,j}(x_i^*) - \rho_{j,i}(x_j^*) = 0 \quad \forall (i,j) \in \mathcal{E} \right).$$

We inductively prove over the tree height that each summand is 0 for our candidate $x^* \in \Omega$.

1. Let i^* be the root. We show that

$$\begin{aligned} & \sum_{j \in \tilde{\mathcal{N}}(i^*)} \min_{x_j \in \Omega_j} f_{i^*,j}(x_{i^*}^*, x_j) - \rho_{i^*,j}(x_{i^*}^*) - \rho_{j,i^*}(x_j) \\ &= -\rho_{i^*} + \sum_{j \in \tilde{\mathcal{N}}(i^*)} \min_{x_j \in \Omega_j} f_{i^*,j}(x_{i^*}^*, x_j) - \rho_{j,i^*}(x_j) \\ & \quad (i_1 \text{ is a root; definition of } \rho_{i^*} \text{ in Theorem 4.1}) \\ &= -\rho_{i^*} + \min_{x_{i^*}^* \in \Omega_{i^*}} \sum_{j \in \tilde{\mathcal{N}}(i^*)} m_{i^*,j}(x_{i^*}^*) \quad (\text{definition of } m_{i^*,j}; \text{ definition of } x_{i^*}^*) \\ &= 0, \quad (\text{by Postulate 3.2; definition of } \rho_{i^*} \text{ in Theorem 4.1}) \end{aligned}$$

Which implies that $x_j \in \Omega_j$ can be selected for all vertices j of depth one, such that the the associated summand is zero, i.e.

$$0 = \min_{x_j \in \Omega_j} f_{i^*,j}(x_{i^*}^*, x_j) - \rho_{i^*,j}(x_{i^*}^*) - \rho_{j,i^*}(x_j) \quad \forall j \in \tilde{\mathcal{N}}(i^*).$$

2. Given the presented method actually selects a successive value $x_j^* \in \Omega_j$ in this way, i.e.

$$0 = f_{i,j}(x_i^*, x_j^*) - \rho_{i,j}(x_i^*) - \rho_{j,i}(x_j^*)$$

for some $(i,j) \in \mathcal{E}$. We know that

$$\begin{aligned} & \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i\}} \min_{x_k \in \Omega_k} f_{j,k}(x_j^*, x_k) - \rho_{j,k}(x_j^*) - \rho_{k,j}(x_k) \\ &= \left(f_{i,j}(x_i^*, x_j^*) - \rho_{i,j}(x_i^*) - \rho_{j,i}(x_j^*) \right) \\ & \quad + \left(\sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i\}} \min_{x_k \in \Omega_k} f_{j,k}(x_j^*, x_k) - \rho_{j,k}(x_j^*) - \rho_{k,j}(x_k) \right) \\ & \quad (\text{assumption; adding 0}) \\ &= \left(\min_{x_i \in \Omega_i} f_{i,j}(x_i, x_j^*) - \rho_{i,j}(x_i) - \rho_{j,i}(x_j^*) \right) \\ & \quad + \left(\sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i\}} \min_{x_k \in \Omega_k} f_{j,k}(x_j^*, x_k) - \rho_{j,k}(x_j^*) - \rho_{k,j}(x_k) \right) \\ & \quad (\text{by } f_{i,j} - \rho_{i,j} - \rho_{j,i} \geq 0) \\ &= -\rho_j + \left(\min_{x_i \in \Omega_i} f_{i,j}(x_i, x_j^*) - \rho_{i,j}(x_i) + \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i\}} m_{j,k}(x_j^*) \right) \\ & \quad (\text{definition of } m_{j,k} \text{ and } \rho_j \text{ in Theorem 4.1}) \\ &= -\rho_j + \min_{x_j \in \Omega_j} \left(\min_{x_i \in \Omega_i} f_{i,j}(x_i, x_j) - \rho_{i,j}(x_i) + \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i\}} m_{j,k}(x_j) \right) \\ & \quad (\text{definition of } x_j^*; \text{swapping minimization}) \end{aligned}$$

$$= 0. \quad (\text{by Postulate 3.2; definition of } \rho_{i^*} \text{ in Theorem 4.1})$$

This implies that all values $x_k \in \Omega_k$ of child vertices k of j , can also be selected, such that their associated summand is zero, i.e.

$$0 = \min_{x_k \in \Omega_k} f_{j,k}(x_j^*, x_k) - \rho_{j,k}(x_j^*) - \rho_{k,j}(x_k) \quad \forall k \in \tilde{\mathcal{N}}(j) \setminus \{i\}.$$

3. The method, in fact, selects successive values such that summands are minimized to be zero, i.e.

$$\begin{aligned} & \arg \min_{x_j \in \Omega_j} f_{i_j,j}(x_{i_j}^*, x_j) + \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i_j\}} m_{j,k}(x_j) \\ &= \arg \min_{x_j \in \Omega_j} f_{i_j,j}(x_{i_j}^*, x_j) - \rho_{i_j,j}(x_{i_j}^*) - \rho_j + \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i_j\}} \min_{x_k \in \Omega_k} f_{j,k}(x_j, x_k) - \rho_{k,j}(x_k) \\ & \quad (\text{definition of } m_{j,k}; \rho_j, \rho_{i_j,j}(x_{i_j}^*) \in \mathbb{R}) \\ &= \arg \min_{x_j \in \Omega_j} f_{i_j,j}(x_{i_j}^*, x_j) - \rho_{i_j,j}(x_{i_j}^*) - \rho_{j,i_j}(x_j) \\ & \quad + \sum_{k \in \tilde{\mathcal{N}}(j) \setminus \{i_j\}} \min_{x_k \in \Omega_k} f_{j,k}(x_j, x_k) - \rho_{j,k}(x_j) - \rho_{k,j}(x_k) \\ & \quad (\text{definition of } \rho_j \text{ in Theorem 4.1}) \\ &\subseteq \arg \min_{x_j \in \Omega_j} f_{i_j,j}(x_{i_j}^*, x_j) - \rho_{i_j,j}(x_{i_j}^*) - \rho_{j,i_j}(x_j). \end{aligned}$$

(previous result; non-negativity of summands)

□

4. Algorithms

We have developed the necessary results to derive several relevant types of optimization algorithms for sums of bivariates in this section. The optimization algorithms that we present differ in the objective function they operate on and their degree of approximating it. For example, the first algorithm we introduce, coordinate descent (CD), operates directly on the sum of bivariates, which is the problem we ultimately aim to solve. All other algorithms, in contrast, operate on some form of dual relaxation of the sum of bivariates that we covered in the previous section.

While we derived insight into tractable optimization problems in the previous section, we will use these results to derive optimization algorithms for more general but similar problem formulations.

4.1. Coordinate Descent for the Sum of Bivariates

As a baseline algorithm for comparison in experiments, we define a simple coordinate descent procedure to heuristically minimize a sum of bivariates. Since we have shown the problem to be NP-hard, we do not expect efficiency nor convergence to the optimal solution, in general.

In contrast to the coordinate optimization approaches presented below, the following coordinate descent algorithm (CD) sequentially optimizes along single coordinates. The reason is that optimizing over neighborhoods of coordinates could yield an objective that is a partial sum of bivariates that does not have tree structure.

Algorithm 1 Coordinate Descent for the Sum of Bivariates (CD)

Input: vertex set $\mathcal{V} := \mathbb{N}_{\leq n}$, $n \in \mathbb{N}$;
 edge set $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$;
 candidate sets $\Omega_i \subset \mathbb{R}$, where $|\Omega_i| \in \mathbb{N}$ and $i \in \mathcal{V}$;
 functions $f_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}$, where $(i, j) \in \mathcal{E}$;
 bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$;
 budget $B \in \mathbb{N}$.

Output: Element $x^* \in \Omega_1 \times \cdots \times \Omega_n$ with “low” function value $\sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j)$.

Initialize: $\tilde{\mathcal{E}} := \{(i, j) \mid (i, j) \in \mathcal{E}\}$; $\tilde{\mathcal{N}}(i) = \{j \mid (i, j) \in \tilde{\mathcal{E}}\}, i \in \mathcal{V}$.

```

1: for all  $i \in (c^{-1}(1), \dots, c^{-1}(n))$  do
2:    $x_i^* \in \arg \min_{x_i \in \Omega_i} \sum_{j \in \tilde{\mathcal{N}}(i)} f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}})$ 
3: end for
4: if  $B \leq t$  then                                     # stopping criterion
5:   return  $x^*$ 
6: end if
7:  $t \leftarrow t + 1$ 
8: go to line 1                                           # repeat the loop

```

4.2. Linear Programming for the Dual Linear Program

Relaxation is the main alternative to an elementary approach that operates directly on the sum of bivariates. As in the previous section, due to strong duality and fewer variables, the dual of the relaxation replaces the relaxation.

The variables of the dual can be interpreted as constant sums of bivariates, such that each bivariate lower-bounds the respective bivariate of the objective function (note the constraint

in [Theorem 4.1](#)). In the case of [Theorem 4.1](#), we introduce the dual from a primal perspective, as it arguably better motivates constant univariate dual variables, namely as a type of local lower bound. Specifically, the problem we present in [Theorem 4.1](#), generalizes the dual relaxation of [Postulate 3.1](#) to non-star graphs, while [Remark 3](#) explains when it is still exact.

Theorem 4.1 (Dual Linear Program). *In the setting of [Definition 1.1](#), let a sum of bivariates $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j} : \Omega \rightarrow \mathbb{R}$ be given, then*

$$\begin{aligned}
 & \begin{cases} \max_{P \in \mathbb{R}^\Omega} \min_{x \in \Omega} P(x) \\ \text{s.t. } P \text{ is separable} \\ P \equiv \sum_{(i,j) \in \mathcal{E}} \tilde{\rho}_{i,j} \\ \tilde{\rho}_{i,j} \in \mathbb{R}^{\Omega_i \times \Omega_j} \\ \tilde{\rho}_{i,j} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \end{cases} & \text{(Separable Local Lower Bound)} \\
 & = \begin{cases} \max_{\rho_{\cdot,\cdot} \in \mathbb{R}^\Omega} \sum_{i \in \mathcal{V}} \min_{x_i \in \Omega_i} \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j}(x_i) \\ \text{s.t. } \rho_{i,j} + \rho_{j,i} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \end{cases} & \text{(Dual Tree Relaxation)} \\
 & = \begin{cases} \max_{\substack{\rho_{\cdot,\cdot} \in \mathbb{R}^\Omega \\ \rho_i \in \mathbb{R}}} \sum_{i \in \mathcal{V}} \rho_i \\ \text{s.t. } \rho_{i,j} + \rho_{j,i} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \\ \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j} \equiv \rho_i \quad \forall i \in \mathcal{V}. \end{cases} & \text{(Dual Linear Program)}
 \end{aligned}$$

Remark 3. By [Theorem 3.3](#) and [Postulate 3.1](#), we know that for star graphs the minimal value of the Dual Linear Program is the minimal value of the sum of bivariates. By [Example 3.5](#), we know that this does not apply beyond trees in general.

Proof. We prove the equality of the first and third term

$$\begin{aligned}
 & \begin{cases} \max_{P \in \mathbb{R}^\Omega} \min_{x \in \Omega} P(x) \\ \text{s.t. } P \text{ is separable} \\ P \equiv \sum_{(i,j) \in \mathcal{E}} \tilde{\rho}_{i,j} \\ \tilde{\rho}_{i,j} \in \mathbb{R}^{\Omega_i \times \Omega_j} \\ \tilde{\rho}_{i,j} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \end{cases} \\
 & = \begin{cases} \max_{P \in \mathbb{R}^\Omega} \min_{x \in \Omega} P(x) \\ \text{s.t. } P \text{ is constant} \\ P \equiv \sum_{(i,j) \in \mathcal{E}} \tilde{\rho}_{i,j} \\ \tilde{\rho}_{i,j} \in \mathbb{R}^{\Omega_i \times \Omega_j} \\ \tilde{\rho}_{i,j} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \end{cases} & \text{(by projection } P \mapsto (x \mapsto \min P)) \\
 & = \begin{cases} \max_{P \in \mathbb{R}^\Omega} \min_{x \in \Omega} P(x) \\ \text{s.t. } P \text{ is constant} \\ P \equiv \sum_{(i,j) \in \mathcal{E}} \rho_{i,j} + \rho_{j,i} \\ \rho_{i,j} \in \mathbb{R}^{\Omega_i}, \rho_{j,i} \in \mathbb{R}^{\Omega_j} \\ \rho_{i,j} + \rho_{j,i} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \end{cases} & \text{(Corollary 2.1)} \\
 & = \begin{cases} \max_{P \in \mathbb{R}^\Omega} \min_{x \in \Omega} \sum_{i \in \mathcal{V}} \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j}(x_i) \\ \text{s.t. } P \text{ is constant} \\ P \equiv \sum_{(i,j) \in \mathcal{E}} \rho_{i,j} + \rho_{j,i} \\ \rho_{i,j} \in \mathbb{R}^{\Omega_i}, \rho_{j,i} \in \mathbb{R}^{\Omega_j} \\ \rho_{i,j} + \rho_{j,i} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \end{cases} & \text{(rewrite objective)}
 \end{aligned}$$

$$\begin{aligned}
&= \begin{cases} \max_{P \in \mathbb{R}^\Omega} \sum_{i \in \mathcal{V}} \min_{x_i \in \Omega_i} \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j}(x_i) \\ \text{s.t. } P \text{ is constant} \\ P \equiv \sum_{(i,j) \in \mathcal{E}} \rho_{i,j} + \rho_{j,i} \\ \rho_{i,j} \in \mathbb{R}^{\Omega_i}, \rho_{j,i} \in \mathbb{R}^{\Omega_j}, \rho_i, \rho_j \in \mathbb{R} \\ \rho_{i,j} + \rho_{j,i} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \\ \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j} \equiv \rho_i \quad \forall i \in \mathcal{V} \end{cases} \quad (P \text{ constant} \iff \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j} \text{ constant, } i \in \mathcal{V}) \\
&= \begin{cases} \max_{\substack{\rho_{\cdot,\cdot} \in \mathbb{R}^\Omega \\ \rho_{\cdot} \in \mathbb{R}}} \sum_{i \in \mathcal{V}} \rho_i \\ \text{s.t. } \rho_{i,j} + \rho_{j,i} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \\ \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j} \equiv \rho_i \quad \forall i \in \mathcal{V}. \end{cases} \quad (\text{eliminate variables and rewrite objective})
\end{aligned}$$

The first equality becomes clear, by replacing $\rho_i, i \in \mathcal{V}$ in the objective and optimizing over non-constant functions $\sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j}$ instead. \square

In algorithm **LPDLP**, we minimize a sum of bivariates by using the dual linear program formulation of [Theorem 4.1](#) to generate a dual solution and then deriving a solution candidate to the sum of bivariates using the method described in [Theorem 3.6](#). The algorithm is proven to be exact for tree-indexed sums of bivariates, as the tree relaxation of [Theorem 3.3](#) is exact, its dual is strong via a similar argument as in [Postulate 3.1](#), and the recovery of a primal solution is exact via [Theorem 3.6](#).

Algorithm 2 Linear Programming for the Dual Linear Program (LPDLP)

Input: vertex set $\mathcal{V} := \mathbb{N}_{\leq n}$, $n \in \mathbb{N}$;

edge set $\mathcal{E} \subseteq \{(i,j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$;

candidate sets $\Omega_i \subset \mathbb{R}$, where $|\Omega_i| \in \mathbb{N}$ and $i \in \mathcal{V}$;

functions $f_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}$, where $(i,j) \in \mathcal{E}$.

bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$;

Output: Element $x^* \in \Omega_1 \times \dots \times \Omega_n$ with “low” function value $\sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j)$.

Initialize: $\tilde{\mathcal{E}} := \{(i,j) \mid (i,j) \in \mathcal{E}\}$; $\tilde{\mathcal{N}}(i) = \{j \mid \{i,j\} \in \tilde{\mathcal{E}}\}, i \in \mathcal{V}$;

$\rho_{i,j}, m_{i,j} := 0 \in \mathbb{R}^{\Omega_i}, \rho_{j,i}, m_{j,i} := 0 \in \mathbb{R}^{\Omega_j}, (i,j) \in \mathcal{E}; \quad \rho_i := 0 \in \mathbb{R}, i \in \mathcal{V}$.

1: solve for # solving the *dual linear program*; see [Theorem 4.1](#)

$$\begin{cases} (\rho_{\cdot,\cdot}^*, \rho_{\cdot}^*) \in \arg \max_{\substack{\rho_{\cdot,\cdot} \in \mathbb{R}^\Omega \\ \rho_{\cdot} \in \mathbb{R}}} \sum_{i \in \mathcal{V}} \rho_i \\ \text{s.t. } \rho_{i,j}(x_i) + \rho_{j,i}(x_j) \leq f_{i,j}(x_i, x_j) \quad \forall x_i \in \Omega_i \quad \forall x_j \in \Omega_j \quad \forall (i,j) \in \mathcal{E} \\ \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j}(x_i) = \rho_i \quad \forall x_i \in \Omega_i \quad \forall i \in \mathcal{V} \end{cases}$$

2: $m_{i,j} \leftarrow \min_{x_j \in \Omega_j} f_{i,j}(\cdot, x_j) - \rho_{j,i}^*(x_j), (i,j) \in \mathcal{E}$ # dual to primal; see [Theorem 3.6](#)

3: $m_{j,i} \leftarrow \min_{x_i \in \Omega_i} f_{i,j}(x_i, \cdot) - \rho_{i,j}^*(x_i), (i,j) \in \mathcal{E}$

4: **for** $i \in (c^{-1}(1), \dots, c^{-1}(n))$ **do**

5: $x_i^* \in \arg \min_{x_i \in \Omega_i} \sum_{\substack{(i,j) \in \tilde{\mathcal{E}} \\ c(j) < c(i)}} f_{i,j}(x_i, x_j) + \sum_{\substack{(i,j) \in \tilde{\mathcal{E}} \\ c(i) < c(j)}} m_{i,j}(x_i)$

6: **end for**

7: **return** x^*

4.3. Block Coordinate Ascent for the Dual Tree Relaxation

Consider restricting the minimization of the dual tree relaxation from [Theorem 4.1](#) to $\rho_{i^*,j} \in \mathbb{R}^{\Omega_i}, j \in \tilde{\mathcal{N}}(i^*)$ for some coordinate $i^* \in \mathcal{V}$. This yields a new problem on what is called a *block coordinate* of the objective function. We have

$$\begin{aligned} & \begin{cases} \max_{\rho_{i^*,\cdot} \in \mathbb{R}^{\Omega_{i^*}}} \sum_{i \in \mathcal{V}} \min_{x_i \in \Omega_i} \sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j}(x_i) \\ \text{s.t.} \quad \rho_{i,j} + \rho_{j,i} \leq f_{i,j} \quad \forall (i,j) \in \mathcal{E} \end{cases} \\ &= \begin{cases} \max_{\rho_{i^*,\cdot} \in \mathbb{R}^{\Omega_{i^*}}} \min_{x_{i^*} \in \Omega_{i^*}} \sum_{j \in \tilde{\mathcal{N}}(i^*)} \rho_{i^*,j}(x_{i^*}) \\ \text{s.t.} \quad \rho_{i^*,j} \leq \min_{x_j \in \Omega_j} f_{i^*,j}(\cdot, x_j) - \rho_{j,i^*}(x_j) \quad \forall j \in \tilde{\mathcal{N}}(i^*). \end{cases} \\ & \quad \text{(possibly transposing } f_{\cdot,\cdot}; \text{ dropping constant terms)} \end{aligned}$$

However, by [Postulate 3.2](#), we know how to solve this problem. Repeatedly replacing variables by their block coordinate optima in some given order yields our version of what is called (*block*) *coordinate ascent* in [BCADTR](#).

Some authors provide counterexamples even for convergence to an optimal solution of the dual linear program of [Theorem 4.1](#) for methods of this class [[DW22](#), Sec. 4.4]. Others make similar claims [[Wer07](#); [PW16](#); [Tou+18](#); [Sav+19](#)]. Evidence for the hypothesis that [BCADTR](#) does not, in general, converge to an optimal solution of the dual linear program can also be observed in [Section 5.1](#). To obtain a solution candidate for the minimization of the sum of bivariates, we apply the method described in [Theorem 3.6](#) once again.

4.4. Block Coordinate Ascent for the Dual Entropy-Regularized Tree Relaxation

This algorithm is inspired by the entropy-based algorithm of [[OP24](#)]. We showed by [Examples 3.1](#) and [3.5](#) that, in case the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is not a tree, translating solutions of the *tree relaxation* presented in [Theorem 3.3](#) to solutions to the minimization of sums of bivariates is NP-hard. However, as previously noted, we have developed a method in [Theorem 3.6](#) that can still be applied heuristically. Further, using linear programming, we can already obtain a dual solution—see [Section 4.2](#) for details. However, one may expect to use the closed-form solution on block coordinates of [Section 3.2](#) to derive a different, possibly a computationally more efficient, solver. We attempted this already in [Section 4.3](#)

We thus propose a block coordinate ascent for the dual of a generalization of the entropy-regularized star relaxation of [Theorem 3.5](#) to arbitrary graphs, termed *entropy-regularized tree relaxation*. This generalization is similarly regularized by an entropy function and approximates tree-structured sums of bivariates.

The main differences to the block coordinate ascent solver for the dual tree relaxation of [Section 4.3](#) lie in

1. the fact that the objective of the following solver only results in an approximation of the dual of the tree relaxation, and
2. that we are able to prove convergence of the solver to a solution with an objective value that is arbitrarily close to the optimal value of the dual tree relaxation.

We present the algorithm in [BCADETR](#). Note that it is not obvious whether better solution candidates of the dual yield better solution candidates to the sum of bivariates, in general. Therefore, similarly as for [BCADTR](#), it is not clear in which iteration the method obtains the best solution to the original problem of minimizing a sum of bivariates.

Algorithm 3 Block Coordinate Ascent for the Dual Tree Relaxation (BCADTR)

Input: vertex set $\mathcal{V} := \mathbb{N}_{\leq n}$, $n \in \mathbb{N}$;
 edge set $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$;
 candidate sets $\Omega_i \subset \mathbb{R}$, where $|\Omega_i| \in \mathbb{N}$ and $i \in \mathcal{V}$;
 functions $f_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}$, where $(i, j) \in \mathcal{E}$;
 bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$;
 weights $w_{i,j;t} \in \mathbb{R}$, $j \in \tilde{\mathcal{N}}(i)$, $i \in \mathcal{V}$ with $\sum_{j \in \tilde{\mathcal{N}}(i)} w_{i,j;t} = 1$, $t \in \mathbb{N}_{\leq B}$;
 budget $B \in \mathbb{N}$.

Output: Element $x^* \in \Omega_1 \times \dots \times \Omega_n$ with “low” function value $\sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j)$.

Initialize: $t := 1 \in \mathbb{N}$; $\tilde{\mathcal{E}} := \{\{i, j\} \mid (i, j) \in \mathcal{E}\}$;
 $\tilde{\mathcal{N}}(i) = \{j \mid \{i, j\} \in \tilde{\mathcal{E}}\}$, $m_i := 0 \in \mathbb{R}^{\Omega_i}$, $i \in \mathcal{V}$;
 $m_{i,j}, \rho_{i,j} := 0 \in \mathbb{R}^{\Omega_i}$, $m_{j,i}, \rho_{j,i} := 0 \in \mathbb{R}^{\Omega_j}$, $(i, j) \in \mathcal{E}$;
 $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$;
 $x^* \in \Omega_1 \times \dots \times \Omega_n$ uniformly at random.

- 1: **for all** $i \in (c^{-1}(1), \dots, c^{-1}(n))$ **do** # one loop in order defined by c
- 2: **for all** $\ell \in \tilde{\mathcal{N}}(i)$ **do** # block coordinate maximization; see [Postulate 3.2](#)
- 3: $m_{i,j}(x_i) \leftarrow \min_{x_j \in \Omega_j} f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}}) - \rho_{j,i}(x_j)$, $\forall x_i \in \Omega_i$
- 4: **end for**
- 5: $m_i \leftarrow \sum_{j \in \tilde{\mathcal{N}}(i)} m_{i,j}$
- 6: $\rho_{i,j} \leftarrow m_{i,j} - w_{i,j;t}(m_i - \min_{x_i \in \Omega_i} m_i(x_i))$, $j \in \tilde{\mathcal{N}}(i)$
- 7: **end for**
- 8: **for** $i \in (c^{-1}(1), \dots, c^{-1}(n))$ **do** # dual to primal; see [Theorem 3.6](#)
- 9: $y_i^* \in \arg \min_{x_i \in \Omega_i} \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(j) < c(i)}} f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}}) + \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(i) < c(j)}} m_{i,j}(x_i)$
- 10: **end for**
- 11: **if** $F(y^*) < F(x^*)$ **then** # keep best solution candidate
- 12: $x^* \leftarrow y^*$
- 13: **end if**
- 14: **if** $B \leq t$ **then** # stopping criterion
- 15: **return** x^*
- 16: **end if**
- 17: $t \leftarrow t + 1$
- 18: **go to line 1** # repeat the loop

Remark 4 (Representation of Dual Entropy-Regularized Tree Relaxation). *In the setting of [Definition 1.1](#), let a sum of bivariates $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j} : \Omega \rightarrow \mathbb{R}$ and $\varepsilon > 0$ be given. Then, by a similar argument as in [Theorem 3.5](#) and [Lemma A.2](#), we have*

$$\begin{aligned}
 & \left(\min_{\substack{\mu_{\cdot, \cdot} \in \mathcal{M}^+(\Omega_{\cdot} \times \Omega_{\cdot}) \\ \mu_{\cdot, \cdot} \in \mathcal{M}_1^+(\Omega_{\cdot})}} \left(\sum_{(i,j) \in \mathcal{E}} \langle f_{i,j}, \mu_{i,j} \rangle - \varepsilon H(\mu_{i,j}) \right) + \varepsilon \sum_{i \in \mathcal{V}} (|\tilde{\mathcal{N}}(i)| - 1) \cdot H(\mu_i) \right) \\
 & \quad \text{s.t.} \quad \mu_i = \mu_{i,j} \circ \pi_i^{-1}, \quad \forall (i, j) \in \mathcal{E} \\
 & \quad \quad \mu_j = \mu_{i,j} \circ \pi_j^{-1}, \quad \forall (i, j) \in \mathcal{E}
 \end{aligned}$$

(entropy-regularized tree relaxation)

$$= \begin{cases} \max_{\rho_{\cdot} \in \mathbb{R}} \max_{\rho_{\cdot, \cdot} \in \mathbb{R}^{\Omega_{\cdot}}} \left(\sum_{i \in \mathcal{V}} \rho_i \right) - \varepsilon \sum_{(i,j) \in \mathcal{E}} \sum_{\substack{x_i \in \Omega_i \\ x_j \in \Omega_j}} \exp \left((\rho_{i,j}(x_i) + \rho_{j,i}(x_j) - f_{i,j}(x_i, x_j)) / \varepsilon - 1 \right) \\ \text{s.t. } \left(\sum_{j \in \tilde{\mathcal{N}}(i)} \rho_{i,j} \right)(x_i) = \rho_i, \forall x_i \in \Omega_i, i \in \mathcal{V}. \end{cases}$$

(dual entropy-regularized tree relaxation)

Remark 5. Due to boundedness of the entropy H , it is clear that for a given sum of bivariates, the entropy-regularized tree relaxation uniformly approximates the tree relaxation of [Theorem 3.3](#). This holds, of course, in particular, if the graph associated with the sum of bivariates is a star graph, as in [Theorem 3.5](#).

Corollary 4.1. The objective value of a block coordinate ascent w.r.t. the blocks $\rho_{i,j}, \rho_i; j \in \tilde{\mathcal{N}}(i)$ for all $i \in \mathcal{V}$ converges to the global maximal value of the dual entropy-regularized tree relaxation of [Remark 4](#).

Proof. Clearly, due to the convexity of \exp , the *dual entropy-regularized tree relaxation* is a sum of concave functions in $\rho_{\cdot} \in \mathbb{R}, \rho_{\cdot, \cdot} \in \mathbb{R}^{\Omega_{\cdot}}$. The objective function is smooth, and the domain of the function is a block-wise product space. Furthermore, by [Theorem 3.5](#), block coordinate-wise minimization yields a unique minimum. Therefore, we can apply [Theorem C.1](#). \square

Algorithm 4 Block Coordinate Ascent for the
Dual Entropy-Regularized Tree Relaxation (BCADETR)

Input: vertex set $\mathcal{V} := \mathbb{N}_{\leq n}$, $n \in \mathbb{N}$;
 edge set $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$;
 candidate sets $\Omega_i \subset \mathbb{R}$, where $|\Omega_i| \in \mathbb{N}$ and $i \in \mathcal{V}$;
 functions $f_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}$, where $(i, j) \in \mathcal{E}$;
 entropy regularization coefficient $\varepsilon \in (0, \infty)$;
 bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$;
 budget $B \in \mathbb{N}$.

Output: Element $x^* \in \Omega_1 \times \dots \times \Omega_n$ with “low” function value $\sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j)$.

Initialize: $t := 1 \in \mathbb{R}$; $\tilde{\mathcal{E}} := \{\{i, j\} \mid (i, j) \in \mathcal{E}\}$;
 $\tilde{\mathcal{N}}(i) = \{j \mid \{i, j\} \in \tilde{\mathcal{E}}\}$,
 $\gamma_i := \varepsilon |\tilde{\mathcal{N}}(i)| + \varepsilon |\tilde{\mathcal{N}}(i)| \log(|\tilde{\mathcal{N}}(i)|/\varepsilon)$, $i \in \mathcal{V}$;
 $\rho_{i,j}, m_{i,j}, \underline{m}_{i,j}, \overline{m}_{i,j} := 0 \in \mathbb{R}^{\Omega_i}$, $\rho_{j,i}, m_{j,i}, \underline{m}_{j,i}, \overline{m}_{j,i} := 0 \in \mathbb{R}^{\Omega_j}$, $(i, j) \in \mathcal{E}$;
 $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$;
 $x^* \in \Omega_1 \times \dots \times \Omega_n$ uniformly at random.

```

1: for all  $i \in (c^{-1}(1), \dots, c^{-1}(n))$  do                                # one loop in order defined by  $c$ 
2:   for all  $\ell \in \tilde{\mathcal{N}}(i)$  do                                                # block coordinate maximization; see Theorem 3.5
3:      $\overline{m}_{i,\ell}(x_i) \leftarrow \text{lse}_{\varepsilon}^{x_{\ell}}(f_{\min\{i,\ell\}, \max\{i,\ell\}}(x_{\min\{i,\ell\}}, x_{\max\{i,\ell\}}) - \rho_{\ell,i}(x_{\ell})), \forall x_i \in \Omega_i$ 
4:      $\underline{m}_{i,\ell}(x_i) \leftarrow -\text{lse}_{\varepsilon}^{x_{\ell}}(\rho_{\ell,i}(x_{\ell}) - f_{\min\{i,\ell\}, \max\{i,\ell\}}(x_{\min\{i,\ell\}}, x_{\max\{i,\ell\}})), \forall x_i \in \Omega_i$ 
5:   end for
6:    $\rho \leftarrow \gamma_i - |\tilde{\mathcal{N}}(i)| \text{lse}_{\varepsilon}^{x_i, \ell} \left( \overline{m}_{i,\ell}(x_i) - \underline{m}_{i,\ell}(x_i) - \frac{1}{|\tilde{\mathcal{N}}(i)|} \sum_{(i,k) \in \tilde{\mathcal{N}}(i)} \overline{m}_{i,k}(x_i) \right)$ 
7:    $\rho_{i,j} \leftarrow \overline{m}_{i,j} + \frac{1}{|\tilde{\mathcal{N}}(i)|} \left( \rho - \sum_{(i,k) \in \tilde{\mathcal{N}}(i)} \overline{m}_{i,k} \right), \forall j \in \tilde{\mathcal{N}}(i)$ 
8: end for
9:  $m_{i,j} \leftarrow \min_{x_j \in \Omega_j} f_{i,j}(\cdot, x_j) - \rho_{j,i}(x_j), (i, j) \in \mathcal{E}$                                 # dual to primal; see Theorem 3.6
10:  $m_{j,i} \leftarrow \min_{x_i \in \Omega_i} f_{i,j}(x_i, \cdot) - \rho_{i,j}(x_i), (i, j) \in \mathcal{E}$ 
11: for  $i \in (c^{-1}(1), \dots, c^{-1}(n))$  do
12:    $y_i^* \in \arg \min_{x_i \in \Omega_i} \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(j) < c(i)}} f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}}) + \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(i) < c(j)}} m_{i,j}(x_i)$ 
13: end for
14: if  $F(y^*) < F(x^*)$  then                                                # keep best solution candidate
15:    $x^* \leftarrow y^*$ 
16: end if
17: if  $B \leq t$  then                                                        # stopping criterion
18:   return  $x^*$ 
19: end if
20:  $t \leftarrow t + 1$ 
21: go to line 1                                                            # repeat the loop

```

5. Experiments

In this section, we compare implementations of the algorithms **CD**, **LPDLP**, **BCADTR**, **TRW-S** and **TRW-S-LEG** for the optimization of several qualitatively distinct and practically relevant objective function classes.

The algorithms **TRW-S** and **TRW-S-LEG** are considered state-of-the-art methods for the optimization of sums of bivariates. The experiments include two versions of **BCADTR**. Namely one, where the weights are sampled uniformly on the simplex for each coordinate and each time step independently, termed “random”, and another version, where the weights are constant for each coordinate and each time step, termed “constant”. The numerically stable implementation of **BCADETR** is considered to be beyond the scope of this work.

The analysis of the value of dual tree relaxation is restricted to **LPDLP** and both versions of **BCADTR** due to their consistent parameterization.

All algorithms are available in the Python programming language and are implemented in a vectorized, sparse as well as in-place manner, whenever beneficial.¹ The implementation of **LPDLP** is based on the *HiGHS solver* for linear programs [HH18]. For all iterative algorithms that operate on the objective, intermediate results are shown.

In addition, *wall-clock time* is chosen for comparing the algorithms as it is a common index.

5.1. Random Sums of Bivariates

First, we construct random sums of bivariates. As usual, we are in the setting of [Definition 1.1](#), and we assume the function values of all bivariates to be independently normal distributed. In addition, the edges of the graph that indexes the sum of bivariates follow the distribution of the $G(n, m)$ -model of Erdős–Rényi random graphs. Precisely, we have

$$\begin{aligned} \Omega_i &:= \{1, \dots, K_i\}, \quad K_i \sim \text{Unif}(\{5, \dots, 15\}) \quad \forall i \in \mathbb{N}_{\leq n}, \\ F &\equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}, \\ f_{i,j}(x_i, x_j) &\sim \mathcal{N}(0, 1) \text{ independently} \quad \forall x_i \in \Omega_i \quad \forall x_j \in \Omega_j \quad \forall (i, j) \in \mathcal{E}, \text{ and} \\ \mathcal{E} &\sim \text{Unif}(\{\mathcal{H} \subseteq \mathcal{J} \mid |\mathcal{H}| = m\}), \text{ where } \mathcal{J} := \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}, m \in \mathbb{N}. \end{aligned}$$

For the experiment, we select the number of arguments of the problem as $n = 100$, the bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$ uniformly at random and the number of edges m , such that the density $m/|\mathcal{J}|$ of the graph $(\mathcal{V}, \mathcal{E})$ approximates the values 10^{-2} , 10^{-1} and $9 \cdot 10^{-1}$. We run the algorithms on 10 independent samples of the objective function $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$. The results are visualized in [Figure 5](#).

An interpretation of the results of the first experiment is that sufficiently dense random sums of bivariates contain structure that is hard to exploit for algorithms based on dual relaxations. In particular, even a simple method, such as (primal) coordinate descent (**CD**) outperforms all other algorithms on sufficiently dense instances. Additionally, random sums of bivariates seem to lie outside the design domain of the legacy tree-reweighted message passing (**TRW-S-LEG**), as it is outperformed by the block-coordinate ascent algorithm for the dual tree relaxation (**BCADTR**).

Solving the dual tree relaxation more precisely does not in general yield better primal candidates, as can be seen by the exact solution to the dual linear program (**LPDLP**). Interestingly, for sufficiently dense random instances, the dual block coordinate ascent algorithm

¹An implementation can be found at <https://github.com/NiMlr/pybiv>.

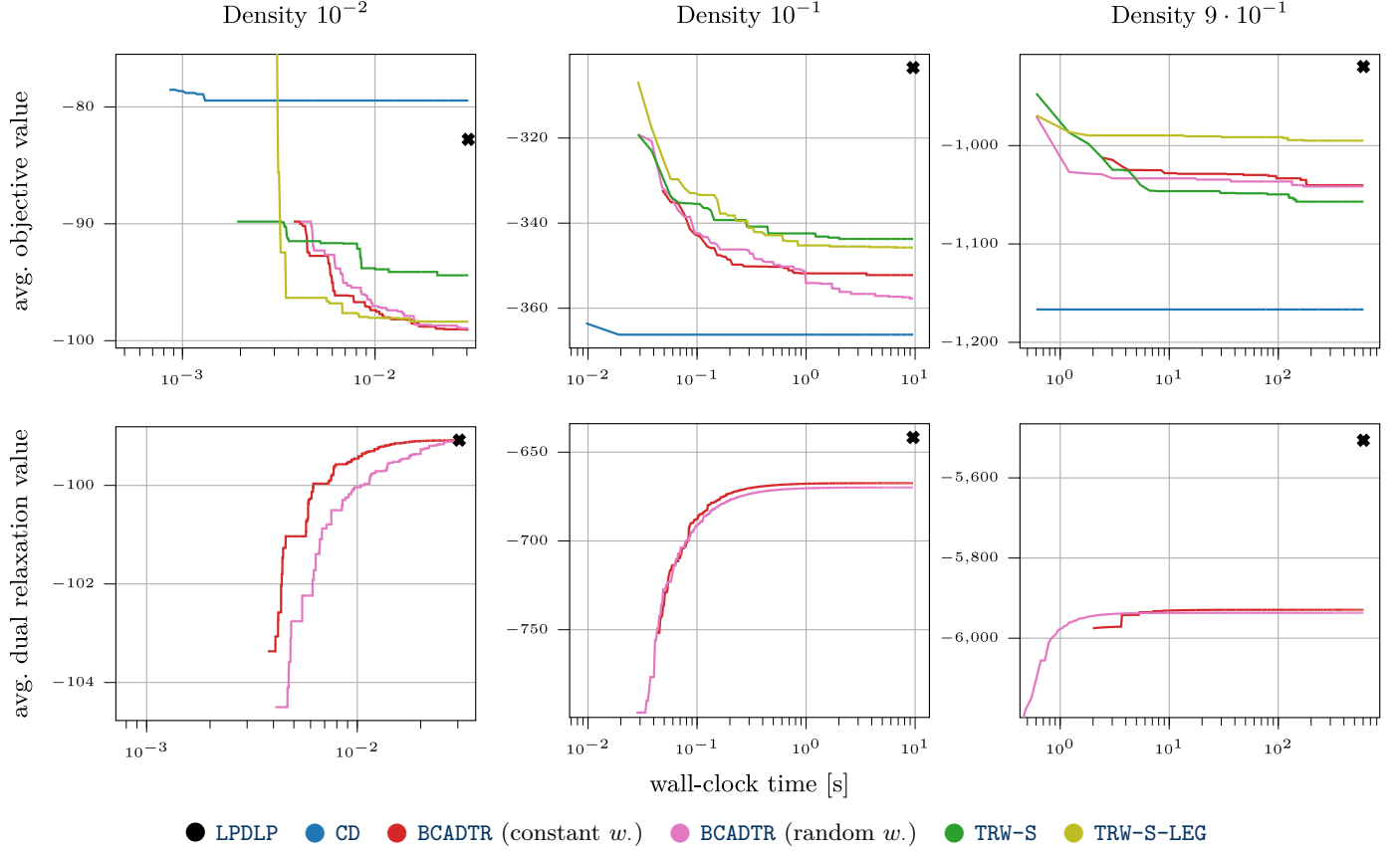


Figure 5: Minimization of random sums of bivariate. Benchmark of several algorithms for the optimization of sums of bivariate with 100 arguments that can each take 5 – 15 values for various densities of the Erdős-Rényi random graph that indexes the sum of bivariate. The plots represent the average of 10 independent runs and of the best value found after given wall-clock time. Primal and dual relaxed values are shown. See Section 5.1 for details.

(BCADTR) does not converge to the maximal value of the dual linear program. It can also be seen that the quality of the primal candidate derived from the solution to the dual linear program deteriorates for dense instances at least as much as the quality of the best candidates of all other algorithms.

5.2. Vertex Coloring

For a given graph $\mathcal{G} = (\mathcal{E}, \mathcal{V})$, vertex coloring asks to determine whether we can assign one of n_c colors for each vertex, such that neighboring vertices are assigned different colors. Mathematically, the goal is to determine the existence of a function $\text{color} : \mathcal{V} \rightarrow \{1, \dots, n_c\}$, $n_c \in \mathbb{N}$, such that $\text{color}(i) \neq \text{color}(j)$ for all edges $(i, j) \in \mathcal{E}$.

The problem of determining whether a n_c -coloring with $n_c = 4$ exists for the vertices of a graph is NP-complete. In the next experiment, we attempt to find a 4-coloring for a given random graph with low density by minimizing the number of color defects between neighboring vertices.

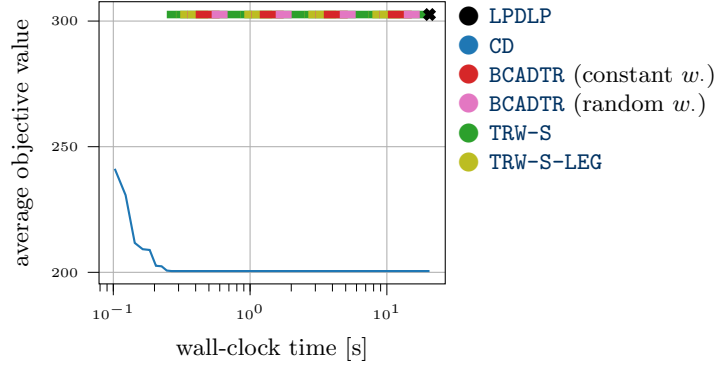


Figure 6: Minimization of the number of defects in vertex coloring. Benchmark of several algorithms for the optimization of sums of bivariates with 10^3 arguments. Each argument that represent one of 4 colors in a vertex coloring of a given Erdős–Rényi random graph with density 10^{-2} that indexes the sum of bivariates. The plots represent the average of 10 independent runs and of the best value found after given wall-clock time. See Section 5.2 for details.

Precisely, we have in the setting of Definition 1.1 that

$$\begin{aligned}\Omega_i &:= \{1, 2, 3, 4\}, \\ F &\equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}, \\ f_{i,j}(x_i, x_j) &= \begin{cases} 1 & \text{if } x_i = x_j \\ 0 & \text{else} \end{cases} \quad \forall x_i \in \Omega_i \quad \forall x_j \in \Omega_j \quad \forall (i,j) \in \mathcal{E}, \text{ and} \\ \mathcal{E} &\sim \text{Unif}(\{\mathcal{H} \subseteq \mathcal{J} \mid |\mathcal{H}| = m\}), \text{ where } \mathcal{J} := \{(i,j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}, m \in \mathbb{N}.\end{aligned}$$

Clearly, if we determine $x \in \Omega$ such that $F(x) = 0$, then x encodes a 4-coloring of the graph \mathcal{G} .

For the experiment, we select the number of arguments of the problem as $n = 1000$, the bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$ uniformly at random, and the number of edges m , such that the density $m/|\mathcal{J}|$ of the graph $(\mathcal{V}, \mathcal{E})$ approximates 10^{-2} . We run the algorithms on 10 independent samples of the objective function $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$. The results are visualized in Figure 6.

An interpretation of the results of the second experiment is that all algorithms that are based on (tree) relaxations can only determine trivial solutions if the structure of the bivariates is sufficiently adverse. This holds even when the density of the graph is comparatively low. Only coordinate descent (CD) progresses by sequentially selecting each vertex color such that it has the least number of defects. However, we must conclude that all considered algorithms are not even heuristically competitive for the problem class of vertex coloring.

5.3. Signal Reconstruction

The challenge of signal reconstruction describes a setting in which a mathematical function, termed *signal*, is superimposed with (random) noise and where the goal is to determine the function without noise or a good approximation of it, termed *reconstruction*. For our purposes, we consider a binary signal $\text{sig} : \mathbb{Z}_n \rightarrow \{0, 1\}$ defined on a discretized circle, the quotient group \mathbb{Z}_n . The noisy version of the signal is given by $\text{sig}_{\text{noisy}} : \mathbb{Z}_n \rightarrow \mathbb{R}$, where $\text{sig}_{\text{noisy}}(x) = \text{sig}(x) + \nu_x$ and $\nu_x \sim \text{Unif}([-1, 1])$ independently for all $x \in \mathbb{Z}_n$.

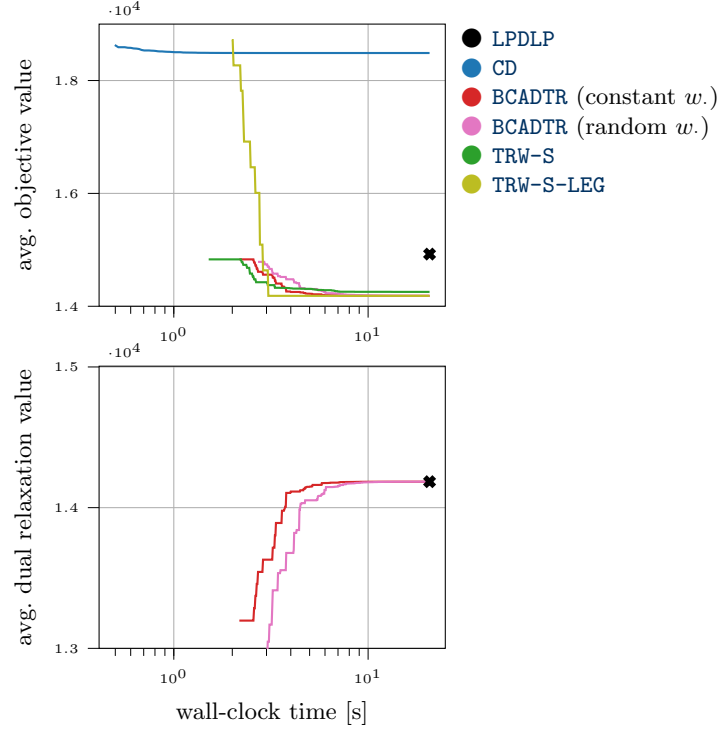


Figure 7: Minimization of the error in a (regularized) signal reconstruction problem.

Benchmark of several algorithms for the optimization of sums of bivariates with $25 \cdot 10^3$ arguments. Each argument represents one of 2 values the reconstruction of a binary signal can take. The graph that indexes the sum of bivariates representing the error has an approximate density of $8 \cdot 10^{-5}$. The plots represent the average of 10 independent runs and of the best value found after given wall-clock time. Primal and dual relaxed values are shown. See Section 5.3 for details.

Given the noisy version of the signal, our reconstruction of the signal is determined by minimizing a heuristic *error function*, which commonly happens to have a representation as a sum of bivariates. The error function penalizes the reconstruction both when deviating from the noisy signal and when neighboring values are different. Precisely, we have in the setting of Definition 1.1 that

$$\begin{aligned}
 \Omega_i &:= \{0, 1\}, \\
 F &\equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}, \\
 f_{i,j}(x_i, x_j) &= |\text{sig}_{\text{noisy}}(x_i) - x_i| \\
 &\quad + (1/2) \cdot \begin{cases} 0 & \text{if } x_i = x_j \\ 1 & \text{else} \end{cases} \quad \forall x_i \in \Omega_i \quad \forall x_j \in \Omega_j \quad \forall (i,j) \in \mathcal{E}, \text{ and} \\
 \mathcal{E} &= \{(x, x+1) \mid x \in \mathbb{Z}_n\}.
 \end{aligned}$$

For the experiment, we select the number of arguments of the problem as $n = 25 \cdot 10^3$, the signal sig is sampled uniformly at random and the bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$ as $c(i) = i$ for all $i \in \mathcal{V}$. The density is approximately $8 \cdot 10^{-5}$. We run the algorithms on 10 independent samples of the objective function $F \equiv \sum_{(i,j) \in \mathcal{E}} f_{i,j}$. The results are visualized in Figure 7.

An interpretation of the results of the third experiment is that all algorithms that are based on (tree) relaxations perform well in reconstructing the signal. The fact that the relaxation gap is rather small and that all such methods outperform **CD** significantly suggests that the performance of the algorithms is non-trivial. Legacy tree-reweighted message passing (**TRW-S-LEG**) works especially well—arguably because signal reconstruction methods are the design domain of this method. This class of problems seems to be the only class in which explicitly solving the dual linear program (**LPDLP**) is competitive. In addition, the block-coordinate ascent methods converge to the maximum of the dual linear program.

6. Conclusion and Future Work

We provided a comprehensive introduction to relaxation theory for the optimization of sums of bivariate functions that resulted in closed-form, linear programming, and block coordinate ascent solution approaches. The work was motivated by a fundamental perspective on both hard as well as easy instances of sums of bivariate functions, and a distributional perspective on the function class using a (no-)free-lunch theorem.

The theory featured insightful results on the limits of relaxation-based approaches, as well as fundamental restrictions on the reconstruction of global functions from their bivariate marginals, which are intimately related to the theory of sums of bivariate functions. A central technique when working with bivariate functions was the design of methods that not only provide the correct result for tree-structured sums of bivariate functions, but are applicable outside this narrow domain. We continued to apply this principle in entropy calculations and regularization.

Promising future work extends the theory of tractable subclasses of the sums of bivariate functions. This could be done, e.g., via a theory of separable lower bounds to bivariate functions as well as via approximation results that extend our analysis. Although suggested by the complexity results and the relaxation theory of this work, our experiments provide further evidence for this claim.

Generally, it is plausible that a theory of sums of bivariate functions shall precede a theory of the minimization thereof. It is also plausible that such generic structures as sums of bivariate functions are useful in other areas of applied mathematics.

Alternatives to relaxation-based optimization algorithms may be found via integer linear programming formulations of sums of bivariate functions.

References

- [Ama72] S.-I. Amari. “Learning Patterns and Pattern Sequences by Self-Organizing Nets of Threshold Elements”. In: *IEEE Transactions on Computers* C-21.11 (1972), pp. 1197–1206.
- [Ber18] D. P. Bertsekas. *Nonlinear Programming: 3rd Edition*. Athena Scientific, 2018.
- [DW22] T. Dlask and T. Werner. “Classes of linear programs solvable by coordinate-wise minimization”. In: *Annals of Mathematics and Artificial Intelligence* 90.7 (2022), pp. 777–807.
- [GJ07] A. Globerson and T. Jaakkola. “Approximate inference using conditional entropy decompositions”. In: *Proceedings of Machine Learning Research*. PMLR, 2007, pp. 131–138.
- [GK02] C. Geiger and C. Kanzow. *Theorie und Numerik restringierter Optimierungsaufgaben*. Springer Berlin Heidelberg, 2002.
- [HH18] Q. Huangfu and J. J. Hall. “Parallelizing the dual revised simplex method”. In: *Mathematical Programming Computation* 10.1 (2018), pp. 119–142.
- [IT04] C. Igel and M. Toussaint. “A No-Free-Lunch Theorem for Non-Uniform Distributions of Target Functions”. In: *Journal of Mathematical Modelling and Algorithms* 3.4 (2004), pp. 313–322.
- [Kap+15] J. H. Kappes et al. “A Comparative Study of Modern Inference Techniques for Structured Discrete Energy Minimization Problems”. In: *International Journal of Computer Vision* 115.2 (2015), pp. 155–184.
- [Kol05] V. Kolmogorov. “Convergent tree-reweighted message passing for energy minimization”. In: *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*. Vol. R5. *Proceedings of Machine Learning Research*. PMLR, 2005, pp. 182–189.
- [Kol14] V. Kolmogorov. “A new look at reweighted message passing”. In: *IEEE transactions on pattern analysis and machine intelligence* 37.5 (2014), pp. 919–930.
- [KPT07] N. Komodakis, N. Paragios, and G. Tziritas. “MRF Optimization via Dual Decomposition: Message-Passing Revisited”. In: *2007 IEEE 11th International Conference on Computer Vision*. 2007, pp. 1–8.
- [KV06] B. Korte and J. Vygen. *Combinatorial Optimization: Theory and Algorithms*. Springer, 2006.
- [LI13] Q. Liu and A. Ihler. “Variational algorithms for marginal MAP”. In: *Journal of Machine Learning Research* 14 (2013), pp. 3165–3200.
- [LS21] J.-H. Lange and P. Swoboda. “Efficient Message Passing for 0–1 ILPs with Binary Decision Diagrams”. In: *Proceedings of the 38th International Conference on Machine Learning*. Vol. 139. *Proceedings of Machine Learning Research*. PMLR, 2021, pp. 6000–6010.
- [LSH16] M. Li, A. Shekhovtsov, and D. Huber. “Complexity of Discrete Energy Minimization Problems”. In: *Computer Vision – ECCV 2016*. 2016, pp. 834–852.
- [Mar52] H. Markowitz. “Portfolio Selection”. In: *The Journal of Finance* 7.1 (1952), pp. 77–91.
- [MG21] N. Müller and T. Glasmachers. “Non-local optimization: imposing structure on optimization problems by relaxation”. In: *Proceedings of the 16th ACM/SIGEVO Conference on Foundations of Genetic Algorithms*. FOGA ’21. Association for Computing Machinery, 2021.

- [OP24] P. Ochs and T. Pock. “Optimization View on DBCA”. In: *Unpublished* (2024).
- [PW15] D. Průša and T. Werner. “Universality of the Local Marginal Polytope”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37.4 (2015), pp. 898–904.
- [PW16] D. Průša and T. Werner. “LP relaxation of the Potts labeling problem is as hard as any linear program”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39.7 (2016), pp. 1469–1475.
- [Sav+19] B. Savchynskyy et al. *Discrete Graphical Models—An Optimization Perspective*. Vol. 11. 3-4. Now Publishers, Inc., 2019, pp. 160–429.
- [Sch76] M. I. Schlesinger. “Syntactic analysis of two-dimensional visual signals in the presence of noise”. In: *Cybernetics* 12.4 (1976), pp. 612–628.
- [SG07] M. I. Schlesinger and V. V. Giginyak. “Solution to structural recognition (MAX,+)-problems by their equivalent transformations”. In: *Part 1* (2007), pp. 3–15.
- [SK75] D. Sherrington and S. Kirkpatrick. “Solvable Model of a Spin-Glass”. In: *Phys. Rev. Lett.* 35 (26 1975), pp. 1792–1796.
- [Sla59] M. Slater. “Lagrange Multipliers Revisited”. In: *Cowles Foundation Discussion Papers*. 80. 1959.
- [SVW01] C. Schumacher, M. D. Vose, and L. D. Whitley. “The No Free Lunch and problem description length”. In: *Proceedings of the 3rd Annual Conference on Genetic and Evolutionary Computation*. GECCO’01. 2001, pp. 565–570.
- [Tou+18] S. Tourani et al. “MPLP++: Fast, parallel dual block-coordinate ascent for dense graphical models”. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018, pp. 251–267.
- [Tou+20] S. Tourani et al. “Taxonomy of Dual Block-Coordinate Ascent Methods for Discrete Energy Minimization”. In: *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*. Vol. 108. Proceedings of Machine Learning Research. PMLR, 2020, pp. 2775–2785.
- [Wer07] T. Werner. “A Linear Programming Approach to Max-Sum Problem: A Review”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29.7 (2007), pp. 1165–1179.
- [WJ08] M. J. Wainwright and M. I. Jordan. “Graphical Models, Exponential Families, and Variational Inference”. In: *Foundations and Trends in Machine Learning* 1.1–2 (2008), pp. 1–305.

A. Additional Proofs

A.1. Lemmata for Theorem 3.5

Lemma A.1. *In the setting of Definition 1.1 and Theorem 3.5, we have that*

$$\min_{\mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})} \left\langle -\varepsilon(|\mathcal{N}(i^*)| - 1) \log \mu_{i^*} + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}, \mu_{i^*} \right\rangle = \min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right)(x_{i^*}),$$

where $\varepsilon > 0$ and $\rho_{i^*, j} \in \mathbb{R}^{\Omega_{i^*}}, (i^*, j) \in \mathcal{N}(i^*)$.

Proof. W.l.o.g., we assume that $|\mathcal{N}(i^*)| - 1 > 0$. Otherwise, the statement is clearly correct. Define

$$\begin{aligned} G(\mu_{i^*}) &:= \left\langle -\varepsilon(|\mathcal{N}(i^*)| - 1) \log \mu_{i^*} + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}, \mu_{i^*} \right\rangle \\ &= \sum_{x_{i^*} \in \Omega_{i^*}} \mu_{i^*}(x_{i^*}) \left(-\varepsilon(|\mathcal{N}(i^*)| - 1) \log \mu_{i^*}(x_{i^*}) + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}(x_{i^*}) \right) \\ &\quad \text{(definition scalar product)} \end{aligned}$$

for all $\mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})$. We have for all $x_{i^*}, y_{i^*} \in \Omega_{i^*}$ and all strictly positive measures $\mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})$ that

$$\begin{aligned} &\frac{\partial}{\partial \mu_{i^*}(y_{i^*})} \frac{\partial}{\partial \mu_{i^*}(x_{i^*})} G(\mu_{i^*}) \\ &= \frac{\partial}{\partial \mu_{i^*}(y_{i^*})} \frac{\partial}{\partial \mu_{i^*}(x_{i^*})} \mu_{i^*}(x_{i^*}) \left(-\varepsilon(|\mathcal{N}(i^*)| - 1) \log \mu_{i^*}(x_{i^*}) + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}(x_{i^*}) \right) \\ &\quad \text{(definition } G) \\ &= \frac{\partial}{\partial \mu_{i^*}(y_{i^*})} -\varepsilon(|\mathcal{N}(i^*)| - 1) (\log \mu_{i^*}(x_{i^*}) + 1) + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}(x_{i^*}) \quad \text{(product rule)} \\ &= \chi_{x_{i^*}=y_{i^*}} \frac{-\varepsilon(|\mathcal{N}(i^*)| - 1)}{\mu_{i^*}(x_{i^*})}. \quad (\chi \text{ denotes indicator function}) \end{aligned}$$

Therefore, the Hessian of G is negative definite and G is strictly concave. This implies that the minimum of G lies on the boundary of $\mathcal{M}_1^+(\Omega_{i^*})$. Which in turn implies that at least one atom $x_{i^*} \in \Omega_{i^*}$ of the minimum has zero measure. Fixing $\mu_{i^*}(x_{i^*}) = 0$ and minimizing G over the remaining atoms yields a problem with the same properties, recursively. Therefore, the minimum of G must be a Dirac measure. Thus, we have

$$\begin{aligned} &\min_{\mu_{i^*} \in \mathcal{M}_1^+(\Omega_{i^*})} \left\langle -\varepsilon(|\mathcal{N}(i^*)| - 1) \log \mu_{i^*} + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}, \mu_{i^*} \right\rangle \\ &= \min_{\mu_{i^*} \in \mathcal{D}(\Omega_{i^*})} \left\langle -\varepsilon(|\mathcal{N}(i^*)| - 1) \log \mu_{i^*} + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}, \mu_{i^*} \right\rangle \\ &\quad \text{(\mathcal{D}(\Omega_{i^*}) denotes the Dirac measures on } \Omega_{i^*}) \\ &= \min_{\mu_{i^*} \in \mathcal{D}(\Omega_{i^*})} \sum_{x_{i^*} \in \Omega_{i^*}} \mu_{i^*}(x_{i^*}) \left(-\varepsilon(|\mathcal{N}(i^*)| - 1) \log \mu_{i^*}(x_{i^*}) + \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}(x_{i^*}) \right) \\ &\quad \text{(definition scalar product)} \\ &= \min_{\mu_{i^*} \in \mathcal{D}(\Omega_{i^*})} \sum_{x_{i^*} \in \Omega_{i^*}} \chi_{\mu_{i^*}(x_{i^*})=1} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}(x_{i^*}) \quad (\chi \text{ denotes indicator function}) \end{aligned}$$

$$= \min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right) (x_{i^*}). \quad \square$$

Lemma A.2. *In the setting of Definition 1.1 and Theorem 3.5, we have for all $(i^*, j) \in \mathcal{N}(i^*)$ that*

$$\begin{aligned} & \min_{\mu_{i^*, j} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_j)} \left\langle \varepsilon \log \mu_{i^*, j} + f_{i^*, j} - \rho_{i^*, j}, \mu_{i^*, j} \right\rangle \\ &= -\varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{x_j \in \Omega_j} \exp \left((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j)) / \varepsilon - 1 \right), \end{aligned}$$

where $\varepsilon > 0$ and $\rho_{i^*, j} \in \mathbb{R}^{\Omega_{i^*}}$.

Proof. We have

$$\begin{aligned} & \min_{\mu_{i^*, j} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_j)} \left\langle \varepsilon \log \mu_{i^*, j} + f_{i^*, j} - \rho_{i^*, j}, \mu_{i^*, j} \right\rangle \\ &= \min_{\mu_{i^*, j} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_j)} \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{x_j \in \Omega_j} \mu_{i^*, j}(x_{i^*}, x_j) \left(\varepsilon \log \mu_{i^*, j}(x_{i^*}, x_j) + f_{i^*, j}(x_{i^*}, x_j) - \rho_{i^*, j}(x_{i^*}) \right) \\ &= \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{x_j \in \Omega_j} \min_{\mu_{i^*, j}(x_{i^*}, x_j) \in [0, \infty)} \mu_{i^*, j}(x_{i^*}, x_j) \left(\varepsilon \log \mu_{i^*, j}(x_{i^*}, x_j) + f_{i^*, j}(x_{i^*}, x_j) - \rho_{i^*, j}(x_{i^*}) \right). \end{aligned}$$

It can be observed that for all $x_{i^*} \in \Omega_{i^*}$ and $x_j \in \Omega_j$, the minimizer of each subproblem will never be at ∞ , since the objective converges to ∞ for $\mu_{i^*, j}(x_{i^*}, x_j) \rightarrow \infty$. Further, at $\mu_{i^*, j}(x_{i^*}, x_j) \downarrow 0$ the objective has the value 0.

Therefore, we check a first order criterion to find critical points with $\mu_{i^*, j}(x_{i^*}, x_j) > 0$. We have for such critical points that

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{d}{d\mu_{i^*, j}(x_{i^*}, x_j)} \mu_{i^*, j}(x_{i^*}, x_j) \left(\varepsilon \log \mu_{i^*, j}(x_{i^*}, x_j) + f_{i^*, j}(x_{i^*}, x_j) - \rho_{i^*, j}(x_{i^*}) \right) \\ &= \varepsilon (\log \mu_{i^*, j}(x_{i^*}, x_j) + 1) + f_{i^*, j}(x_{i^*}, x_j) - \rho_{i^*, j}(x_{i^*}) \quad (\text{product rule}) \\ &\iff \mu_{i^*, j}(x_{i^*}, x_j) = \exp \left((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j)) / \varepsilon - 1 \right). \end{aligned}$$

Plugging the critical point into the objective, we see that the associated objective value is

$$\begin{aligned} & \exp \left((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j)) / \varepsilon - 1 \right) \left(\varepsilon \log \exp \left((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j)) / \varepsilon - 1 \right) \right. \\ & \quad \left. + f_{i^*, j}(x_{i^*}, x_j) - \rho_{i^*, j}(x_{i^*}) \right) \\ &= -\exp \left((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j)) / \varepsilon - 1 \right) \varepsilon. \end{aligned}$$

The critical point, which has a negative function value, must be a minimizer, as at both boundaries the function has non-negative values and there are no further critical points.

Plugging this representation of the minimal value into the initial equation, we get

$$\begin{aligned} & \min_{\mu_{i^*, j} \in \mathcal{M}^+(\Omega_{i^*} \times \Omega_j)} \left\langle \varepsilon \log \mu_{i^*, j} + f_{i^*, j} - \rho_{i^*, j}, \mu_{i^*, j} \right\rangle \\ &= -\varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{x_j \in \Omega_j} \exp \left((\rho_{i^*, j}(x_{i^*}) - f_{i^*, j}(x_{i^*}, x_j)) / \varepsilon - 1 \right). \quad \square \end{aligned}$$

Lemma A.3. *In the setting of Definition 1.1 and Theorem 3.5, we have*

$$\begin{aligned}
& \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\left(\min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right) (x_{i^*}) \right) \right. \\
& \quad \left. - \varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \right). \\
& = \begin{cases} \rho - |\mathcal{N}(i^*)| \\ \text{where} \\ \rho_{i^*, j}(x_{i^*}) = \text{lse}_\varepsilon(f_{i^*, j}(x_{i^*}, \cdot)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon(f_{i^*, k}(x_{i^*}, \cdot)) \right), \\ \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*) \\ \rho = \varepsilon |\mathcal{N}(i^*)| + \varepsilon |\mathcal{N}(i^*)| \log(|\mathcal{N}(i^*)|/\varepsilon) \\ - |\mathcal{N}(i^*)| \text{lse}_\varepsilon^{x_{i^*}, j} \left(\text{lse}_\varepsilon^{x_j}(f_{i^*, j}(x_{i^*}, x_j)) + \text{lse}_\varepsilon^{x_j}(-f_{i^*, j}(x_{i^*}, x_j)) \right. \\ \quad \left. - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon^{x_k}(f_{i^*, k}(x_{i^*}, x_k)) \right). \end{cases}
\end{aligned}$$

Proof. We have

$$\begin{aligned}
& \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \left(\left(\min_{x_{i^*} \in \Omega_{i^*}} \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right) (x_{i^*}) \right) \right. \\
& \quad \left. - \varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \right) \\
& = \begin{cases} \max_{\rho \in \mathbb{R}} \max_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \rho - \varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\ \text{s.t. } \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right) (x_{i^*}) = \rho, \forall x_{i^*} \in \Omega_{i^*} \end{cases} \\
& \hspace{25em} (\text{separability \& monotonicity}) \\
& = \max_{\rho \in \mathbb{R}} \rho - \varepsilon \cdot \begin{cases} \min_{\rho_{i^*, \cdot} \in \mathbb{R}^{\Omega_{i^*}}} \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\ \text{s.t. } \left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j} \right) (x_{i^*}) = \rho, \forall x_{i^*} \in \Omega_{i^*} \end{cases} \\
& \hspace{25em} (\text{separability}) \\
& = \begin{cases} \max_{\rho \in \mathbb{R}} \rho - \varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\ \text{where } \rho_{i^*, j}(x_{i^*}) = \text{lse}_\varepsilon(f_{i^*, j}(x_{i^*}, \cdot)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon(f_{i^*, k}(x_{i^*}, \cdot)) \right), \\ \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*) \end{cases} \\
& \hspace{25em} (\text{see “inner problem” below})
\end{aligned}$$

$$\begin{aligned}
& \left\{ \begin{aligned} & \rho - \varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\ & \text{where} \\ & \rho_{i^*, j}(x_{i^*}) = \text{lse}_\varepsilon(f_{i^*, j}(x_{i^*}, \cdot)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon(f_{i^*, k}(x_{i^*}, \cdot)) \right), \\ & \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*) \\ & \rho = \varepsilon |\mathcal{N}(i^*)| + \varepsilon |\mathcal{N}(i^*)| \log(|\mathcal{N}(i^*)|/\varepsilon) \\ & - |\mathcal{N}(i^*)| \text{lse}_\varepsilon^{x_{i^*}, j} \left(\text{lse}_\varepsilon^{x_j}(f_{i^*, j}(x_{i^*}, x_j)) + \text{lse}_\varepsilon^{x_j}(-f_{i^*, j}(x_{i^*}, x_j)) \right. \\ & \quad \left. - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon^{x_k}(f_{i^*, k}(x_{i^*}, x_k)) \right) \end{aligned} \right. \quad (\text{see “outer problem” below}) \\
& = \left\{ \begin{aligned} & \rho - |\mathcal{N}(i^*)| \\ & \text{where} \\ & \rho_{i^*, j}(x_{i^*}) = \text{lse}_\varepsilon(f_{i^*, j}(x_{i^*}, \cdot)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon(f_{i^*, k}(x_{i^*}, \cdot)) \right), \\ & \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*) \\ & \rho = \varepsilon |\mathcal{N}(i^*)| + \varepsilon |\mathcal{N}(i^*)| \log(|\mathcal{N}(i^*)|/\varepsilon) \\ & - |\mathcal{N}(i^*)| \text{lse}_\varepsilon^{x_{i^*}, j} \left(\text{lse}_\varepsilon^{x_j}(f_{i^*, j}(x_{i^*}, x_j)) + \text{lse}_\varepsilon^{x_j}(-f_{i^*, j}(x_{i^*}, x_j)) \right. \\ & \quad \left. - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon^{x_k}(f_{i^*, k}(x_{i^*}, x_k)) \right). \end{aligned} \right. \quad (\text{see “simplify maximal value” below})
\end{aligned}$$

Simplify maximal value: In the above setting, we have

$$\begin{aligned}
& \rho - \varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\
& = \rho - \varepsilon \exp \frac{1}{\varepsilon} \log \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp \left(\left(\rho_{i^*, j}(x_{i^*}) + \text{lse}_\varepsilon^{x_j}(-f_{i^*, j}(x_{i^*}, x_j) - \varepsilon) \right) / \varepsilon \right) \\
& \quad (\text{definition lse}_\varepsilon; \text{properties of exp; exp log} \equiv \text{id}) \\
& = \rho - \varepsilon \exp \left(\text{lse}_\varepsilon^{x_{i^*}, j} \left(\rho_{i^*, j}(x_{i^*}) + \text{lse}_\varepsilon^{x_j}(-f_{i^*, j}(x_{i^*}, x_j)) \right) / \varepsilon - 1 \right) \\
& \quad (\text{definition lse}_\varepsilon; \text{properties of lse}_\varepsilon) \\
& = \rho - \varepsilon \exp \left(\text{lse}_\varepsilon^{x_{i^*}, j} \left(\text{lse}_\varepsilon^{x_j}(f_{i^*, j}(x_{i^*}, x_j)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon^{x_k}(f_{i^*, k}(x_{i^*}, x_k)) \right) \right. \right. \\
& \quad \left. \left. + \text{lse}_\varepsilon^{x_j}(-f_{i^*, j}(x_{i^*}, x_j)) \right) / \varepsilon - 1 \right) \\
& \quad (\text{representation of } \rho_{i^*, j}(x_{i^*}); \text{renaming variables}) \\
& = \rho - \varepsilon \exp \left(\text{lse}_\varepsilon^{x_{i^*}, j} \left(\text{lse}_\varepsilon^{x_j}(f_{i^*, j}(x_{i^*}, x_j)) - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon^{x_k}(f_{i^*, k}(x_{i^*}, x_k)) \right. \right. \\
& \quad \left. \left. + \text{lse}_\varepsilon^{x_j}(-f_{i^*, j}(x_{i^*}, x_j)) \right) / \varepsilon + \frac{\rho}{\varepsilon |\mathcal{N}(i^*)|} - 1 \right) \\
& \quad (\text{constant addition to lse}_\varepsilon)
\end{aligned}$$

$$\begin{aligned}
&= \rho - \varepsilon \exp \left(1 + \log (|\mathcal{N}(i^*)|/\varepsilon) - 1 \right) && \text{(representation of } \rho) \\
&= \rho - |\mathcal{N}(i^*)|. && \text{(simplifying)}
\end{aligned}$$

Outer problem: Since the outer problem is unbounded from below at infinity, a necessary condition for the maximizers of the outer problem can be found using a first order criterion, i.e. in critical points. Define

$$\begin{aligned}
R(\rho) &:= \rho - \varepsilon \cdot \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\
&\quad \text{where } \rho_{i^*, j}(x_{i^*}) = \underbrace{\text{lse}_\varepsilon(f_{i^*, j}(x_{i^*}, \cdot)) - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon(f_{i^*, k}(x_{i^*}, \cdot)) + \frac{\rho}{|\mathcal{N}(i^*)|}}_{=: \gamma(x_{i^*}, j)}.
\end{aligned}$$

We have

$$\begin{aligned}
0 &= \frac{d}{d\rho} R(\rho) \\
&= 1 - \frac{\varepsilon}{|\mathcal{N}(i^*)|} \cdot \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\
&= 1 - \frac{\varepsilon}{|\mathcal{N}(i^*)|} \cdot \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\gamma(x_{i^*}, j)/\varepsilon + \frac{\rho}{\varepsilon |\mathcal{N}(i^*)|}) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\
&\hspace{15em} \text{(definition of } \rho_{i^*, \cdot}) \\
&\iff \\
1 &= \exp\left(\frac{\rho}{\varepsilon |\mathcal{N}(i^*)|}\right) \frac{\varepsilon}{|\mathcal{N}(i^*)|} \cdot \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\gamma(x_{i^*}, j)/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\
&\hspace{15em} \text{(rearranging)} \\
&\iff \rho = \varepsilon |\mathcal{N}(i^*)| \log \left(\frac{|\mathcal{N}(i^*)|}{\varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\gamma(x_{i^*}, j)/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1)} \right) \\
&\hspace{15em} \text{(rearranging)} \\
&= \varepsilon |\mathcal{N}(i^*)| \log \left(\frac{|\mathcal{N}(i^*)|}{\varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\gamma(x_{i^*}, j)/\varepsilon) \exp(\text{lse}_\varepsilon(-f_{i^*, j}(x_{i^*}, \cdot) - \varepsilon)/\varepsilon)} \right) \\
&\hspace{15em} \text{(definition lse}_\varepsilon) \\
&= \varepsilon |\mathcal{N}(i^*)| \log \left(\frac{|\mathcal{N}(i^*)|}{\varepsilon \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp\left(\frac{1}{\varepsilon}(\gamma(x_{i^*}, j) + \text{lse}_\varepsilon(-f_{i^*, j}(x_{i^*}, \cdot) - \varepsilon))\right)} \right) \\
&\hspace{15em} \text{(property of exp)} \\
&= \varepsilon |\mathcal{N}(i^*)| \log (|\mathcal{N}(i^*)|/\varepsilon) \\
&\quad - \varepsilon |\mathcal{N}(i^*)| \log \left(\sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp\left(\frac{1}{\varepsilon}(\gamma(x_{i^*}, j) + \text{lse}_\varepsilon(-f_{i^*, j}(x_{i^*}, \cdot) - \varepsilon))\right) \right) \\
&\hspace{15em} \text{(property of log)} \\
&= \varepsilon |\mathcal{N}(i^*)| \log (|\mathcal{N}(i^*)|/\varepsilon) - |\mathcal{N}(i^*)| \text{lse}_\varepsilon^{x_{i^*}, j} \left(\gamma(x_{i^*}, j) + \text{lse}_\varepsilon^{x_j} (-f_{i^*, j}(x_{i^*}, x_j) - \varepsilon) \right) \\
&\hspace{5em} \text{(definition of lse}_\varepsilon; \text{superscript of lse}_\varepsilon \text{ denotes the variable being maximized)}
\end{aligned}$$

$$\begin{aligned}
&= \varepsilon |\mathcal{N}(i^*)| + \varepsilon |\mathcal{N}(i^*)| \log (|\mathcal{N}(i^*)|/\varepsilon) \\
&\quad - |\mathcal{N}(i^*)| \text{lse}_{\varepsilon}^{x_{i^*}, j} \left(\gamma(x_{i^*}, j) + \text{lse}_{\varepsilon}^{x_j} (-f_{i^*, j}(x_{i^*}, x_j)) \right) \\
&\hspace{15em} (\text{constant addition to } \text{lse}_{\varepsilon}) \\
&= \varepsilon |\mathcal{N}(i^*)| + \varepsilon |\mathcal{N}(i^*)| \log (|\mathcal{N}(i^*)|/\varepsilon) \\
&\quad - |\mathcal{N}(i^*)| \text{lse}_{\varepsilon}^{x_{i^*}, j} \left(\text{lse}_{\varepsilon}^{x_j} (f_{i^*, j}(x_{i^*}, x_j)) \right. \\
&\quad \left. - \frac{1}{|\mathcal{N}(i^*)|} \sum_{(i^*, k) \in \mathcal{N}(i^*)} \text{lse}_{\varepsilon}^{x_k} (f_{i^*, k}(x_{i^*}, x_k)) + \text{lse}_{\varepsilon}^{x_j} (-f_{i^*, j}(x_{i^*}, x_j)) \right) \\
&\hspace{10em} (\text{definition } \gamma(x_{i^*}, j); \text{ superscript of } \text{lse}_{\varepsilon} \text{ denotes the variable being maximized})
\end{aligned}$$

Inner problem: Since the inner problem is unbounded from above at infinity, a necessary condition for the minimizers of the inner problem can be found using Lagrange multipliers, i.e. in critical points of the function

$$\begin{aligned}
G(\lambda, \rho_{i^*, \cdot}(\cdot)) &:= \sum_{x_{i^*} \in \Omega_{i^*}} \sum_{(i^*, j) \in \mathcal{N}(i^*)} \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \\
&\quad - \sum_{x_{i^*} \in \Omega_{i^*}} \lambda_{x_{i^*}} \left(\left(\sum_{(i^*, j) \in \mathcal{N}(i^*)} \rho_{i^*, j}(x_{i^*}) \right) (x_{i^*}) - \rho \right),
\end{aligned}$$

where $\lambda \in \mathbb{R}^{\Omega_{i^*}}$. We have for a critical point that

$$\begin{aligned}
0 &= \frac{\partial}{\partial \rho_{i^*, j}(x_{i^*})} G(\lambda, \rho_{i^*, \cdot}(\cdot)) \\
&= (1/\varepsilon) \exp(\rho_{i^*, j}(x_{i^*})/\varepsilon) \left(\sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \right) - \lambda_{x_{i^*}} \\
&\iff \varepsilon \log \varepsilon \lambda_{x_{i^*}} \\
&= \rho_{i^*, j}(x_{i^*}) + \varepsilon \log \left(\sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \right), \quad \forall (i^*, j) \in \mathcal{N}(i^*), \quad \forall x_{i^*} \in \Omega_{i^*}. \\
&\hspace{15em} (\text{applying the function } \varepsilon \log(\varepsilon(\cdot + \lambda_{x_{i^*}})))
\end{aligned}$$

From the constraint of the inner problem, we can find the minimizer, as

$$\begin{aligned}
& \left(\sum_{(i^*, k) \in \mathcal{N}(i^*)} \rho_{i^*, k}(x_{i^*}) \right) (x_{i^*}) = \rho, \quad \forall x_{i^*} \in \Omega_{i^*} \\
&\iff \sum_{(i^*, k) \in \mathcal{N}(i^*)} \varepsilon \log \varepsilon \lambda_{x_{i^*}} - \varepsilon \log \left(\sum_{x_k \in \Omega_k} \exp(-f_{i^*, k}(x_{i^*}, x_k)/\varepsilon - 1) \right) = \rho, \quad \forall x_{i^*} \in \Omega_{i^*} \\
&\hspace{15em} (\text{necessary condition on critical points}) \\
&\iff |\mathcal{N}(i^*)| \varepsilon \log \varepsilon \lambda_{x_{i^*}} - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \varepsilon \log \left(\sum_{x_k \in \Omega_k} \exp(-f_{i^*, k}(x_{i^*}, x_k)/\varepsilon - 1) \right) = \rho, \quad \forall x_{i^*} \in \Omega_{i^*} \\
&\hspace{15em} (\text{constant summation}) \\
&\iff |\mathcal{N}(i^*)| \left(\rho_{i^*, j}(x_{i^*}) + \varepsilon \log \left(\sum_{x_j \in \Omega_j} \exp(-f_{i^*, j}(x_{i^*}, x_j)/\varepsilon - 1) \right) \right) \\
&\quad - \sum_{(i^*, k) \in \mathcal{N}(i^*)} \varepsilon \log \left(\sum_{x_k \in \Omega_k} \exp(-f_{i^*, k}(x_{i^*}, x_k)/\varepsilon - 1) \right) = \rho, \quad \forall x_{i^*} \in \Omega_{i^*}, \quad (i^*, j) \in \mathcal{N}(i^*) \\
&\hspace{15em} (\text{necessary condition on critical points})
\end{aligned}$$

$$\begin{aligned}
\Longleftrightarrow \rho_{i^*,j}(x_{i^*}) &= \frac{1}{|\mathcal{N}(i^*)|} \left(\rho + \sum_{(i^*,k) \in \mathcal{N}(i^*)} \varepsilon \log \left(\sum_{x_k \in \Omega_k} \exp(-f_{i^*,k}(x_{i^*}, x_k)/\varepsilon - 1) \right) \right) \\
&\quad - \varepsilon \log \left(\sum_{x_j \in \Omega_j} \exp(-f_{i^*,j}(x_{i^*}, x_j)/\varepsilon - 1) \right), \quad \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*) \\
&\hspace{15em} \text{(rearranging)}
\end{aligned}$$

$$\begin{aligned}
&\Longleftrightarrow \rho_{i^*,j}(x_{i^*}) \\
&= \frac{1}{|\mathcal{N}(i^*)|} \left(\rho + \varepsilon \sum_{(i^*,k) \in \mathcal{N}(i^*)} \log \left(\frac{\sum_{x_j \in \Omega_j} \exp(f_{i^*,j}(x_{i^*}, x_j)/\varepsilon - 1)}{\sum_{x_k \in \Omega_k} \exp(f_{i^*,k}(x_{i^*}, x_k)/\varepsilon - 1)} \right) \right), \quad \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*). \\
&\hspace{15em} \text{(simplifying)}
\end{aligned}$$

$$\begin{aligned}
&\Longleftrightarrow \rho_{i^*,j}(x_{i^*}) \\
&= \text{lse}_\varepsilon(f_{i^*,j}(x_{i^*}, \cdot)) + \frac{1}{|\mathcal{N}(i^*)|} \left(\rho - \sum_{(i^*,k) \in \mathcal{N}(i^*)} \text{lse}_\varepsilon(f_{i^*,k}(x_{i^*}, \cdot)) \right), \quad \forall x_{i^*} \in \Omega_{i^*}, (i^*, j) \in \mathcal{N}(i^*). \\
&\hspace{15em} \text{(simplifying)}
\end{aligned}$$

□

B. Additional Code

B.1. Verifying that $F \circ \tilde{p}$ of [Example 3.4](#) is not a sum of bivariates

```

import numpy as np
from sympy.matrices import Matrix
import itertools
from permutation import Permutation

# drop prefix and add zeros to binary string
def fill(bin_num, n):
    len_bin_num = len(bin_num)-2
    return '0'*(n-len_bin_num)+bin_num[2:]

# evaluate the function given a binary string and its dimension
def F(bin_str_ori):
    n = len(bin_str_ori)
    num = int(bin_str_ori, 2)
    pnum = perm[num]
    bin_str = fill(bin(pnum), n)

    out = 0

    for pair in itertools.combinations(list(range(n)), 2):
        if (bin_str[pair[0]] == bin_str[pair[1]]) & (bin_str[pair[0]] == '1'):
            out += 1
    return out

if __name__ == "__main__":
    # smallest possible dimension for a counterexample
    # with binary domain
    n = 6

    # found by perm = np.random.permutation(2**6)
    perm = np.array([8, 52, 62, 30, 43, 27, 1, 36, 18, 58, 53, 54, 33, 22, 39,
                    49, 23, 10, 11, 57, 29, 6, 25, 32, 47, 16, 55, 56, 26, 0,
                    19, 15, 7, 63, 4, 61, 21, 14, 51, 9, 60, 28, 17, 2, 50, 3,
                    38, 12, 40, 35, 42, 5, 44, 37, 41, 59, 31, 34, 46, 48, 24,
                    45, 20, 13])

    # to cycle notation; external package indexes by 1
    p = Permutation(*list(perm+1))
    p = p.to_cycles()
    print("The permutation is: \n", [tuple(np.array(cycle)-1) for cycle in p])

    # maps parameters of bivariates to the sum of bivariates
    A = np.zeros((2**n, 4*(n*(n-1)//2)), dtype=int)
    # stores function values
    B = np.zeros(2**n, dtype=int)

```

```

# assemble system of eqs
# iterate eq number
for i in range(2*n):
    # iterate bivariate functions
    edge_ind = -1
    # input args of bivariate function
    bin_str = fill(bin(i), n)
    for pair in itertools.combinations(list(range(n)), 2):
        # index of bivariate function
        edge_ind += 1
        inp_args = bin_str[pair[0]] + bin_str[pair[1]]
        # variable number
        var_num = 4*edge_ind + int(inp_args, 2)

        # set term
        A[i, var_num] = 1

    B[i] = F(fill(bin(i), n))

# check system using gauss elimination
A = Matrix(A)
B = Matrix(B)
try:
    A.gauss_jordan_solve(B)
    print("Linear system has a solution", "\nThe composition with the\
permutation results in a sum of bivariates")
except ValueError as msg:
    print(msg, "\nThe composition with the permutation does not\
result in a sum of bivariates")

```


B.2. Pseudocode for TRW-S

Algorithm 5 Legacy Sequential Tree-Reweighted Message Passing (TRW-S-LEG) by [Kol05]

Input: vertex set $\mathcal{V} := \mathbb{N}_{\leq n}$, $n \in \mathbb{N}$;
 edge set $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$;
 candidate sets $\Omega_i \subset \mathbb{R}$, where $|\Omega_i| \in \mathbb{N}$ and $i \in \mathcal{V}$;
 functions $f_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}$, where $(i, j) \in \mathcal{E}$;
 bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$;
 budget $B \in \mathbb{N}$.

Output: Element $x^* \in \Omega_1 \times \dots \times \Omega_n$ with “low” function value $\sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j)$.

Initialize: $t, \delta := 1, 0 \in \mathbb{R}$; $\tilde{\mathcal{E}} := \{(i, j) \mid (i, j) \in \mathcal{E}\}$;
 $m_{i,j} := 0 \in \mathbb{R}^{\Omega_i}$, $m_{j,i} := 0 \in \mathbb{R}^{\Omega_j}$, $(i, j) \in \mathcal{E}$; $m_i := 0 \in \mathbb{R}^{\Omega_i}$, $i \in \mathcal{V}$;
 $p_i := \max \{|\{(i, j) \in \tilde{\mathcal{E}} : c(i) < c(j)\}|, |\{(i, j) \in \tilde{\mathcal{E}} : c(j) < c(i)\}|\}$, $i \in \mathcal{V}$;
 $\gamma_{i,j} := 1/p_i$, $\gamma_{j,i} := 1/p_j$, $(i, j) \in \mathcal{E}$.

- 1: **for all** $i \in (c^{-1}(1), \dots, c^{-1}(n))$ **do** # loop vertices in order defined by c
- 2: $m_i \leftarrow \sum_{\{i,j\} \in \tilde{\mathcal{E}}} m_{i,j}$
- 3: $\delta \leftarrow \min_{x_i \in \Omega_i} m_i(x_i)$
- 4: $m_i(x_i) \leftarrow m_i(x_i) - \delta$, $\forall x_i \in \Omega_i$
- 5: **for all** $\{i, j\} \in \tilde{\mathcal{E}}$ with $c(i) < c(j)$ **do** # restrict to “increasing” edges
- 6: $m_{j,i}(x_j) \leftarrow \min_{x_i \in \Omega_i} \gamma_{i,j} m_i(x_i) - m_{i,j}(x_i) + f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}})$, $\forall x_j \in \Omega_j$
- 7: $\delta \leftarrow \min_{x_j \in \Omega_j} m_{j,i}(x_j)$
- 8: $m_{j,i}(x_j) \leftarrow m_{j,i}(x_j) - \delta$, $\forall x_j \in \Omega_j$
- 9: **end for**
- 10: **end for**
- 11: **for** $i \in (c^{-1}(1), \dots, c^{-1}(n))$ **do** # determine the solution candidate
- 12: $y_i^* \in \arg \min_{x_i \in \Omega_i} \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(j) < c(i)}} f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}}) + \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(i) < c(j)}} m_{i,j}(x_i)$
- 13: **end for**
- 14: **if** $F(y^*) < F(x^*)$ **then** # keep best solution candidate
- 15: $x^* \leftarrow y^*$
- 16: **end if**
- 17: **if** $B \leq t$ **then** # stopping criterion
- 18: **return** x^*
- 19: **end if**
- 20: $t \leftarrow t + 1$ # run again in reversed order
- 21: $c(i) \leftarrow |\mathcal{V}| + 1 - i$, $i \in \mathcal{V}$
- 22: **go to line 1**

Algorithm 6 Sequential Tree-Reweighted Message Passing (TRW-S) by [Tou+20]

Input: vertex set $\mathcal{V} := \mathbb{N}_{\leq n}$, $n \in \mathbb{N}$;
 edge set $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i < j\}$;
 candidate sets $\Omega_i \subset \mathbb{R}$, where $|\Omega_i| \in \mathbb{N}$ and $i \in \mathcal{V}$;
 functions $f_{i,j} : \Omega_i \times \Omega_j \rightarrow \mathbb{R}$, where $(i, j) \in \mathcal{E}$;
 bijection $c : \mathcal{V} \rightarrow \mathbb{N}_{\leq n}$;
 budget $B \in \mathbb{N}$.

Output: Element $x^* \in \Omega_1 \times \dots \times \Omega_n$ with “low” function value $\sum_{(i,j) \in \mathcal{E}} f_{i,j}(x_i, x_j)$.

Initialize: $t := 1 \in \mathbb{N}$; $\tilde{\mathcal{E}} := \{(i, j) \mid (i, j) \in \mathcal{E}\}$;
 $m_{i,j}, \rho_{i,j} := 0 \in \mathbb{R}^{\Omega_i}, m_{j,i}, \rho_{j,i} := 0 \in \mathbb{R}^{\Omega_j}, (i, j) \in \mathcal{E}$; $m_i := 0 \in \mathbb{R}^{\Omega_i}, i \in \mathcal{V}$;
 $p_i := \max \{|\{(i, j) \in \tilde{\mathcal{E}} : c(i) < c(j)\}|, |\{(i, j) \in \tilde{\mathcal{E}} : c(j) < c(i)\}|\}, i \in \mathcal{V}$;
 $\gamma_{i,j} := 1/p_i, \gamma_{j,i} := 1/p_j, (i, j) \in \mathcal{E}$.

- 1: **for all** $i \in (c^{-1}(1), \dots, c^{-1}(n))$ **do** # loop vertices in order defined by c
- 2: $m_i \leftarrow \sum_{\{i,j\} \in \tilde{\mathcal{E}}} m_{i,j}$
- 3: **for all** $\{i, j\} \in \tilde{\mathcal{E}}$ with $c(i) < c(j)$ **do** # restrict to “increasing” edges
- 4: $m_{i,j}(x_i) \leftarrow \min_{x_j \in \Omega_j} f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}}) - \rho_{j,i}(x_j), \forall x_i \in \Omega_i$
- 5: $\rho_{i,j} \leftarrow m_{i,j} - \gamma_{i,j} m_i$
- 6: **end for**
- 7: **end for**
- 8: **for** $i \in (c^{-1}(1), \dots, c^{-1}(n))$ **do** # determine the solution candidate
- 9: $y_i^* \in \arg \min_{x_i \in \Omega_i} \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(j) < c(i)}} f_{\min\{i,j\}, \max\{i,j\}}(x_{\min\{i,j\}}, x_{\max\{i,j\}}) + \sum_{\substack{\{i,j\} \in \tilde{\mathcal{E}} \\ c(i) < c(j)}} m_{i,j}(x_i)$
- 10: **end for**
- 11: **if** $F(y^*) < F(x^*)$ **then** # keep best solution candidate
- 12: $x^* \leftarrow y^*$
- 13: **end if**
- 14: **if** $B \leq t$ **then** # stopping criterion
- 15: **return** x^*
- 16: **end if**
- 17: $t \leftarrow t + 1$ # run again in reversed order
- 18: $c(i) \leftarrow |\mathcal{V}| + 1 - i, i \in \mathcal{V}$
- 19: **go to line 1**

C. Referenced Results

Theorem C.1 (Convergence of Block Coordinate Descent [Ber18, p. 324 Prop. 3.7.1]). *Let*

- *a set $X = X_1 \times X_2 \times \cdots \times X_m$ be given, where $m, n_i \in \mathbb{N}$ and $X_i \subseteq \mathbb{R}^{n_i}$ is closed and convex for all $i \in \mathbb{N}_{\leq m}$,*
- *the function $f : X \rightarrow \mathbb{R}$ be continuously differentiable, and*
- *the function $\xi \in X_i \mapsto f(x_1, \dots, x_{i-1}, \xi, x_{i+1}, \dots, x_m)$ attain a unique minimum $\bar{\xi}$ for all $x \in X$ and all $i \in \mathbb{N}_{\leq m}$, and be monotonically non-increasing in the interval from x_i to $\bar{\xi}$.*

Then, every limit point of the sequence $(x_k)_{k \in \mathbb{N}}$ in X , where $x_0 \in X$, and where for all $k \in \mathbb{N}$, we have

$$x_i^{k+1} \in \arg \min_{\xi \in X_i} f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, \xi, x_{i+1}^k, \dots, x_m^k), \quad i = 1, \dots, m,$$

is a stationary point of f .

Theorem C.2 (Strong Lagrangian Duality [Sla59] & [GK02, p. 322 Thm. 6.13]). *Consider*

- *the convex functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g_i : \mathbb{R}^n \rightarrow \mathbb{R}, i \in \mathbb{N}_{\leq m}$, as well as*
- *the affine functions $h_i : \mathbb{R}^n \rightarrow \mathbb{R}, i \in \mathbb{N}_{\leq p}$, where $n, m, p \in \mathbb{N}$, and*
- *the non-empty convex set $X \subseteq \mathbb{R}^n$.*

Then, if there exists \hat{x} in the relative interior of X and, such that,

$$g_i(\hat{x}) < 0 \quad \text{for } i = 1, \dots, m \quad \text{and} \quad h(\hat{x}) = 0,$$

we have the equality

$$\begin{cases} \inf_{x \in X} f(x) \\ \text{s.t. } g(x) \leq 0 \\ h(x) = 0 \end{cases} = \begin{cases} \sup_{\lambda \in \mathbb{R}^m} \inf_{x \in X} f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x) \\ \text{s.t. } \lambda \geq 0, \end{cases}$$

if the left-hand side of it is finite.

Theorem C.3 (Consistency of Global Relaxations [MG21, Corollary 3.1.1]). *Given*

- *a family of probability measures $\{\mathbb{P}_\theta : \theta \in \Theta\}$ on a measure space (Ω, \mathcal{A}) ,*
- *an optimization problem $\min_{x \in \Omega} f(x)$, where $f : \Omega \rightarrow \mathbb{R}$ is \mathcal{A} -continuous at its unique global minimum $x^* \in \Omega$,*
- *$f \in \mathcal{L}^1(\Omega, \mathcal{A}, \mathbb{P}_\theta)$ for all $\theta \in \Theta$, and*
- *that for all $\varepsilon, \gamma > 0$, we have*

$$\int_{\Omega - U_\gamma(x^*)} \max\{|f|, 1\} d\mathbb{P}_\theta < \varepsilon,$$

where U_γ denotes the γ -ball centered at x .

Then, if a minimum θ^ of*

$$\theta \in \Theta \mapsto \int_{\Omega} f d\mathbb{P}_\theta$$

exists, it is unique and $\mathbb{P}_{\theta^} = \delta_{x^*}$, i.e. the Dirac measure at x^* .*