BOUNDS ON PERFECT NODE CLASSIFICATION: A CONVEX GRAPH CLUSTERING PERSPECTIVE

Anonymous authorsPaper under double-blind review

ABSTRACT

We study the problem of transductive node classification in graphs where communities align with both node features and labels. We propose a novel convex optimization framework that integrates node-specific information (features and labels) into graph clustering via low-rank matrix estimation. Our analysis reveals a bidirectional interaction between graph structure and node information: not only can features aid clustering, but graph structure can also enhance node classification. In particular, we prove that incorporating suitable node information enables perfect recovery of communities under milder conditions than required by graph clustering alone. To make the framework practical, we develop efficient algorithmic solutions and validate our theory with experiments demonstrating the predicted improvements.

1 Introduction

Transductive node classification is a fundamental problem in machine learning on graphs: given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where nodes have features and a subset $\mathcal{V}_{\mathrm{train}}$ is labeled, the goal is to predict the labels of the remaining nodes $\mathcal{V}_{\mathrm{test}}$. This setting is ubiquitous in real-world applications such as citation networks, social networks, and recommender systems (Bhagat et al., 2011; Zhu et al., 2003). A central idea of modern graph-based learning is that exploiting graph structure, in addition to node features, can significantly improve classification accuracy. This idea has driven the development of graph neural networks (GNNs) and related methods (Kipf & Welling, 2017; Hamilton et al., 2017; Veličković et al., 2018).

However, despite strong empirical evidence supporting the utility of graph information, there is little formal understanding of *when and why* it actually helps: no existing work establishes rigorous conditions under which incorporating graph structure provably improves classification performance over feature-only methods.

In this paper, we address this fundamental question by formulating transductive node classification as a convex optimization problem, which enables us to derive rigorous theoretical guarantees about when graph structure provably helps classification. We focus on homophilic graphs (McPherson et al., 2001), where nodes of the same class exhibit higher connectivity than nodes of different classes, a property prevalent in many real-world data, including citation networks (Sen et al., 2008; Namata et al., 2012), social networks (Hamilton et al., 2017), and commercial networks (McAuley et al., 2015). In homophilic settings, the graph topology naturally exhibits community structure that correlates with node features and labels.

Our key insight is that under homophily, graph clustering and node classification are fundamentally complementary: graph clustering exploits the graph topology to reveal community structure, while classification leverages features and partial labels. This observation motivates us to develop a unified framework that integrates both tasks. We build on convex methods for graph clustering (Korlakai Vinayak et al., 2014; Chen et al., 2014; Li et al., 2021), which offer theoretical tractability and are closely related to well-understood spectral techniques (Belkin & Niyogi, 2001; Ng et al., 2001; Hajek et al., 2016). In particular, our approach leverages the framework of atomic norms, which generalize the nuclear norm and allow us to define atoms that jointly capture graph structure and node-specific information. This perspective enables us to extend convex graph clustering formulations based on low-rank positive semidefinite representations allowing them to incorporate node features and labels within a single convex optimization framework.

Our contributions are as follows:

- We introduce a novel convex optimization formulation for transductive node classification that
 generalizes existing graph clustering methods to incorporate node-specific information, including
 features and partial labels. The framework is inspired by principles of multimodal learning and
 provides a principled way to combine structural and attribute information.
- We prove that, under suitable structural assumptions, our framework achieves perfect label recovery, i.e., exact classification of all nodes. To the best of our knowledge, this provides the first rigorous theoretical result demonstrating that node features and graph structure can provably interact to improve node classification.
- We develop CADO, a scalable alternating conditional gradient algorithm for solving our optimization problem with a fixed number of atoms. Each step of the algorithm reduces to tractable subproblems with closed-form solutions, making the method both efficient and practical. Through experiments, we show that CADO solves the proposed optimization and validate our theoretical findings.

To the best of our knowledge, no prior work has formulated node classification through the lens of convex graph clustering. Our framework integrates ideas from convex clustering and graph clustering, while directly incorporating both node features and graph structure into a single formulation.

2 Related Work

Convex clustering. Clustering groups data points based solely on their features (Xu & Wunsch, 2005; Jaeger & Banks, 2023; Shalev-Shwartz & Ben-David, 2014; Soltanolkotabi & Candes, 2012). Classical k-means is NP-hard (Aloise et al., 2009), and Lloyd's algorithm is prone to local minima and sensitive to initialization. Convex clustering addresses these issues by adding a sum-of-norms (SON) regularizer and formulating a convex relaxation of k-means (Hocking et al., 2011; Lindsten et al., 2011; Panahi et al., 2017; Tan & Witten, 2015; Sun et al., 2021). Recovery guarantees have been established for two clusters (Zhu et al., 2014), k clusters (Panahi et al., 2017), and weighted variants (Sun et al., 2021), with further contributions by Chiquet et al. (2017); Chi & Steinerberger (2019). However, these works remain tied to k-means. Our framework allows general loss functions and extending recovery guarantees to this broader setting and to node classification.

Graph clustering. Graph clustering (or community detection) seeks to identify clusters using only graph structure (Schaeffer, 2007; Abbe, 2018; Li et al., 2021). The stochastic block model (Holland et al., 1983) is the standard framework, with exact recovery thresholds established for two clusters (Abbe et al., 2015), general k (Wu et al., 2015), and partially-observed graphs (Chen et al., 2014; Korlakai Vinayak et al., 2014). Our work goes beyond SBM by providing recovery guarantees for deterministic graphs and Erdős-Rényi random graphs, thereby broadening the scope of classical graph clustering theory.

Covariate-assisted graph clustering. Many modern graphs include node covariates, motivating methods that leverage both graph struture and node-level information. Heuristic approaches aggregate the adjacency matrix with a Gram or kernel matrix from covariates (Binkiewicz et al., 2017; Yan & Sarkar, 2021; Hu & Wang, 2024; Chunaev, 2020). The contextual SBM (CSBM) (Deshpande et al., 2018) provides a statistical model for this setting, and recent theory shows that covariates can improve graph clustering (Braun et al., 2022; Dreveton et al., 2023; Yang & Fountoulakis, 2023; Braun & Sugiyama, 2024). These works, however, rely on Gaussian or exponential covariates combined within the CSBM and do not capture explicitly the synergy between graph structure and covariates. Unlike these works, our framework establishes a two-way interaction: not only can features support clustering, but graph structure can also enhance node classification. We prove recovery guarantees for this bidirectional setting under static graphs, arbitrary features, and any convex loss functions.

3 Node Classification via Atomic Norms

In this section, we propose a graph clustering formulation that naturally extends to incorporate node-specific information.

We consider a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where each node $v \in \mathcal{V}$ has feature vector x_v . A subset $\mathcal{V}_{\text{train}} \subset \mathcal{V}$ of the nodes has (possibly noisy) labels, and the goal is to predict the labels for the remaining nodes

 $\mathcal{V}_{\text{test}}$, referred to as test nodes. We denote by y_v the true label of node v and its noisy label by \tilde{y}_v . For convenience, we denote by z_v the information associated with node $v \in \mathcal{V}$, where $z_v = x_v$ for $v \in \mathcal{V}_{\text{test}}$ and $z_v = (x_v, \tilde{y}_v)$ for $v \in \mathcal{V}_{\text{train}}$.

Our starting point is convex optimization techniques for graph clustering, which seek a low-rank positive semi-definite (PSD) approximation of the adjacency matrix A (Korlakai Vinayak et al., 2014),

$$\min_{L \in \mathcal{B}} \|L - A\|_1 + \mu_0 \|L\|_*, \tag{1}$$

where $\|\cdot\|_1$ is the ℓ_1 -norm, $\|\cdot\|_*$ the *nuclear norm* (sum of the singular values), \mathcal{B} denotes the set of symmetric matrices with entries in [0,1] and ones on the diagonal, and $\mu_0 \geq 0$ is a regularization parameter.

Since the nuclear norm is a special case of the atomic norm (Bhaskar et al., 2013) (see Appendix A for details), we can equivalently replace $\|L\|_*$ with the atomic norm $\|L\|_{\mathcal{A}}$, where the atomic set is defined as $\mathcal{A} \coloneqq \{ee^\top : \|e\|_2 = 1\}$, where each atom is a rank-one matrix formed from a unit vector. Using the atomic norm, (1) can be rewritten as

$$\min_{\boldsymbol{L}\in\mathcal{B}} \|\boldsymbol{L} - \boldsymbol{A}\|_1 + \mu_0 \|\boldsymbol{L}\|_{\mathcal{A}}. \tag{2}$$

Noting that $\|L - A\|_1$ can be expressed as $-\langle \bar{A}, L \rangle$, (2) can be rewritten as

$$\min_{\boldsymbol{L}\in\mathcal{B}, \{\lambda_{i}\geq 0, \, \boldsymbol{e}_{i}\in\mathbb{S}^{n-1}\}_{i}} - \langle \bar{\boldsymbol{A}}, \boldsymbol{L} \rangle + \mu_{0} \sum_{i} \lambda_{i}$$
s.t.
$$\boldsymbol{L} = \sum_{i} \lambda_{i} \boldsymbol{e}_{i} \boldsymbol{e}_{i}^{\top},$$
(3)

where A is the polarized adjacency matrix with entries 1 if $\{u,v\} \in \mathcal{E}$ and -1 otherwise. The summations are over the infinite elements as the atomic set contains infinite number of atoms, and this holds for the rest of paper whenever no upper bound is specified.

This optimization serves as the foundation of our framework. Specifically, if \mathcal{G} has r well-separated clusters, the solution of (3) has rank r and admits the decomposition

$$L = E \Lambda E^{\top}, \quad \Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_r), \quad E = [e_1, \dots, e_r],$$

The v-th row of E, ϵ_v represents a low-rank embedding of node $v \in \mathcal{V}$. Under some mild conditions on the graph (see Section 4), these embeddings become one-hot vectors that perfectly indicate cluster membership, enabling exact cluster recovery.

Incorporating node-specific information. We extend the convex graph clustering formulation in (3) to incorporate node-specific information. We associate each node v with a model $\theta_v \in \Theta$ and a loss function $f_v(\mathbf{z}_v; \boldsymbol{\theta})$ measuring how well θ_v fits \mathbf{z}_v . A concrete example will be given in Section 6.1. Based on these components, we formulate the following optimization problem,

$$\min_{\substack{\boldsymbol{L} \in \mathcal{B}, \{\boldsymbol{\theta}_{v} \in \Theta\}_{v \in \mathcal{V}}, \\ \{\lambda_{i} \geq 0, \boldsymbol{e}_{i} \in \mathbb{S}^{n-1}, \boldsymbol{\theta}_{i} \in \Theta\}_{i}}} - \langle \bar{\boldsymbol{A}}, \boldsymbol{L} \rangle + \mu \sum_{v \in \mathcal{V}} f_{v}(\boldsymbol{z}_{v}; \boldsymbol{\theta}_{v})$$

$$\text{s.t.} \quad \boldsymbol{L} = \sum_{i} \lambda_{i} \boldsymbol{e}_{i} \boldsymbol{e}_{i}^{\top}, \quad \boldsymbol{\theta}_{v} = \sum_{i} \lambda_{i} \epsilon_{i,v}^{2} \boldsymbol{\theta}_{i}, \quad \forall v \in \mathcal{V}$$

where $\epsilon_{i,v}$ is the v^{th} element of e_i , and $\theta_i \in \Theta$ represents the model associated with cluster i. The additional term in the objective function of (4) plays the role of an empirical risk, but with each z_v evaluated against a node-specific model θ_v . The coupling constraint $\theta_v = \sum_i \lambda_i \epsilon_{i,v}^2 \theta_i$ ensures these models are not completely independent. In particular, when the embedding vectors $\epsilon_v = (\epsilon_{i,v})_i$ coincide with the one-hot vectors of the corresponding clusters, we have $\theta_v = \theta_i$ for all v in cluster i, so classification benefits directly from the low-rank clustering mechanism.

Similar to (2)–(3), (4) can be expressed as a regularized atomic norm problem. To this end, we define the joint variable $U := (L, \{\theta_v\}_{v \in \mathcal{V}})$, the atoms $a_i := (e_i e_i^{\top}, \{\epsilon_{i,v}^2 \theta_i\}_{v \in \mathcal{V}})$, and the atomic set

$$\mathcal{A} := \left\{ \left(\boldsymbol{e}\boldsymbol{e}^{\top}, \{\epsilon_v^2\boldsymbol{\theta}\}_{v \in \mathcal{V}}\right) \;\middle|\; \boldsymbol{e} = (\epsilon_v)_v, \; \|\boldsymbol{e}\|_2 = 1, \; \boldsymbol{\theta} \in \Theta \right\}.$$

¹This follows from the fact that for $a_{ij} \in \{0,1\}$ and $l_{ij} \in [0,1]$, we have $|l_{ij} - a_{ij}| = -l_{ij}\bar{a}_{ij} + a_{ij}$.

With this notation, problem (4) can be rewritten compactly as

$$\min_{\boldsymbol{U} \in \mathcal{U}, \{\lambda_i \ge 0, \, \boldsymbol{a}_i\}_i} \quad \phi(\boldsymbol{U}) + \mu_0 \sum_i \lambda_i \\
\text{s.t.} \quad \boldsymbol{U} = \sum_i \lambda_i \boldsymbol{a}_i,$$

where $\phi(U) := -\langle \bar{A}, L \rangle + \mu \sum_{v \in \mathcal{V}} f_v(z_v; \theta_v)$ and $\mathcal{U} := \mathcal{B} \times \Theta^{|\mathcal{V}|}$ is the feasible set of all variables $(L, \{\theta_v\}_{v \in \mathcal{V}})$ with $L \in \mathcal{B}$ and $\theta_v \in \Theta$.

Regularization by Sum of Norms. The formulation in (5) admits an arbitrary number of atoms, each potentially corresponding to a distinct cluster. This flexibility risks over-parameterization: the node-specific term in ϕ may favor assigning each node to its own cluster. Although the graph-based term in ϕ can counteract this, the resulting balance does not capture the intended synergy between topology and node-specific information.

To address this, we adopt the well-known *sum-of-norms (SON)* regularization (Lindsten et al., 2011; Panahi et al., 2017) defined as

$$oldsymbol{R}(oldsymbol{U}) := \sum_{u < v} \|oldsymbol{ heta}_u - oldsymbol{ heta}_v\|$$
 .

This penalty encourages node-specific models to coincide, promoting shared representations within clusters. Incorporating it into (5) yields

$$\min_{\boldsymbol{U} \in \mathcal{U}, \{\lambda_i \ge 0, \, \boldsymbol{a}_i\}_i} \quad \phi(\boldsymbol{U}) + \mu_0 \sum_i \lambda_i + \mu_1 R(\boldsymbol{U})$$
s.t.
$$\boldsymbol{U} = \sum_i \lambda_i \boldsymbol{a}_i$$
(6)

This regularization enforces that nodes in the same cluster share identical models, thereby aligning the clustering and classification components.

Equation (6) is our final optimization framework: it extends convex graph clustering by incorporating node-specific information through SON regularization. In the remainder of the paper, we analyze this framework from two complementary perspectives. First, we establish conditions under which it achieves perfect recovery of the underlying clusters. Second, we investigate the computational complexity of solving the problem and introduce CADO, an efficient algorithm for computing its solution.

4 PERFECT RECOVERY

In this section, we establish conditions under which the global solution of (6) achieves perfect recovery of the underlying clusters and, consequently, the class labels. Our analysis proceeds in three steps. First, we characterize the structure of the *ideal solution* that corresponds to perfect recovery. Next, we formalize the probabilistic model governing the graph and node features. Based on this model, we derive explicit conditions on the parameters under which the solution of (6) coincides with the ideal one, thereby achieving perfect recovery. Finally, we present our main theorem, showing that perfect recovery can be achieved from graph structure alone, from node features alone, or from their combination. Crucially, the joint setting admits strictly milder conditions than either source individually, thereby demonstrating the synergistic effect of integrating graph and feature information.

4.1 IDEAL SOLUTION

Consider a population \mathcal{V} of size $|\mathcal{V}| = n$, partitioned into K clusters $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K$, where cluster \mathcal{C}_i contains n_i nodes. In the ideal scenario, the solution consists of K atoms (one per cluster) denoted by $\{(\lambda_i^*, \epsilon_{i,v}^* \boldsymbol{\theta}_i^*)\}_{i=1}^K$, with

$$\lambda_i^* = n_i, \qquad \epsilon_{i,v}^* = \begin{cases} \frac{1}{\sqrt{n_i}}, & v \in \mathcal{C}_i, \\ 0, & \text{otherwise.} \end{cases}$$
 (7)

This yields the ideal partition matrix L^* defined by $L_{uv}^* = 1$ if $u, v \in C_i$ for some i, and 0 otherwise.

Moreover, each cluster \mathcal{C}_i is associated with a single model θ_i^* . These cluster-level models arise from the following characteristic optimization problem, obtained by substituting the ideal coefficients from (7) into (6):

$$\{\boldsymbol{\theta}_i^*\}_i = \min_{\{\boldsymbol{\theta}_i\}_i \in \Theta} \mu \sum_{i=1}^K F_i(\boldsymbol{\theta}_i) + \gamma \sum_{i < j} n_j \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|, \tag{8}$$

where $F_i(\theta) := \sum_{v \in \mathcal{C}_i} f_v(\boldsymbol{z}; \boldsymbol{\theta})$ is the aggregate loss function of cluster \mathcal{C}_i , and $\gamma \coloneqq \mu_1/\mu$. The solutions of (8), referred to as *biased centroids*, satisfy the optimality condition $\mathbf{0} \in \boldsymbol{\omega}_i + \partial I_{\Theta}(\boldsymbol{\theta}_i^*)$, where $\partial I_{\mathcal{U}}(\cdot)$ is the subdifferential of the indicator function $I_{\mathcal{U}}$ of \mathcal{U} , and

$$\boldsymbol{\omega}_i \coloneqq \nabla F_i(\boldsymbol{\theta}_i^*) + \gamma \sum_{j < i} n_j \frac{\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*}{\|\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*\|}.$$
 (9)

4.2 OUR MODEL

Our optimization framework takes as input both the graph structure and the node features. In what follows, we formalize the statistical assumptions underlying each of these components and introduce the definitions that will be used in our recovery analysis.

4.2.1 GRAPH

Let $N_{ij} := \sum_{u \in \mathcal{C}_i, v \in \mathcal{C}_j} A_{uv}$ be the connectivity of the different partitions. Naturally, we are interested in the case where N_{ii} is sufficiently larger than N_{ij} for $i \neq j$. For a node v, let $n_{v,j}$ be the number of edges from v to nodes in cluster \mathcal{C}_j , $n_{v,j} := \sum\limits_{u \in \mathcal{C}_j} A_{uv}$. We take $n_{v,j}^+ = n_j - n_{v,j}$ if $v \notin \mathcal{C}_j$, and $n_{v,j}^+ = n_{v,j}$ if $v \notin \mathcal{C}_j$. We also define $N_{i,j}^+ = \sum\limits_{v \in \mathcal{C}_i} n_{v,j}^+$ and $\rho_{i,j}^+ = N_{i,j}^+/n_i n_j$.

$$v \in \mathcal{C}_j$$
, and $n_{v,j}^+ = n_{v,j}$ if $v \notin \mathcal{C}_j$. We also define $N_{i,j}^+ = \sum_{v \in \mathcal{C}_i} n_{v,j}^+$ and $\rho_{i,j}^+ = N_{i,j}^+/n_i n_j$.

These quantities capture the level of misconnections between clusters and vanish in the ideal case of perfectly separated clusters. For a fixed $\delta > 0$, we formalize and bound the level of cluster separability using the following assumptions.

Assumption 4.1 (δ -Homogeneity). For any node $v \in \mathcal{C}_i$ and any cluster \mathcal{C}_i , it holds that

$$\left| \frac{n_{v,j}^+}{n_j} - \rho_{ij}^+ \right| \le \delta. \tag{10}$$

Note that from Assumption 4.1, the maximum number of mis-connections in each block (i,j), denoted by d_{ij}^{\max} , is bounded as $d_{ij}^{\max} \leq (\rho_{ij}^+ + \delta) n_j$.

Assumption 4.2 (δ -Visibility). For any two clusters C_i , C_j , it holds that $\rho_{ij}^{+} < \frac{1}{2} - \delta$.

4.2.2 Node-Specific Vectors

In convex clustering, perfect recovery requires hat inter-cluster separation dominates intra-cluster variability: loosely, the minimum distance between clusters must be significantly larger than the maximum cluster diameter While prior works focus on the k-means loss function, we consider general convex loss functions, which necessitate more general definitions of inter-cluster distance (separability) and cluster diameter. In the following assumptions, we formalize these two measures.

Assumption 4.3 (*R*-separability). The biased centroids θ_i^* are distinct. Moreover, for every $\theta \in \Theta$, the relation

Figure 1: R-separability.

$$\langle \boldsymbol{\omega}_i, \boldsymbol{\theta} - \boldsymbol{\theta}_i^* \rangle \leq R$$

holds for at most one index i. In addition, the feasible set Θ is assumed to be bounded, i.e.,

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}_i^*\| \leq \ell, \quad \forall i$$

This assumption guarantees that centroids corresponding to different clusters are well-separated, as depicted in Figure 1.

Assumption 4.4 (Gradient variability). For all $u, v \in C_i$, the gradient variability is bounded by $\|\nabla f_v(\theta_i^*) - \nabla f_u(\theta_i^*)\| \le \rho$.

This assumption bounds the variability of node-specific information within clusters, i.e., it quantifies the diameter of the clusters. It ensures that within each cluster, the local losses behave similarly around the corresponding centroid θ_i^* (each node aligns closely with its associated centroid). Furthermore, the parameter ρ can be interpreted as an indirect measure of data noise: larger values of ρ indicate higher within-cluster variability, meaning the data are noisier and less homogeneous.

4.3 Main result

Here, we present our main theoretical results.

Theorem 4.5. Consider a fixed graph \mathcal{G} with n nodes partitioned into K classes, where class i contains n_i nodes with $n_i \sim n_j$ for all i, j. Let ρ_{ij}^+ denote the relative misconnection rate between classes, and define

$$\rho^{+} := \max_{i,j} \rho_{ij}^{+}, \qquad p_{\max} = \rho^{+} + \delta, \qquad a = \max_{i,j} a_{ij}.$$

Fix $\gamma = \mu_1/\mu$, and suppose Assumptions 4.1–4.4 hold with fixed δ and R, and that the feasible set Θ is bounded. Then:

1. **Graph-structure-only regime** ($\mu \rightarrow 0$). Perfect recovery is achievable if

$$p_{\max} \le \frac{1}{2 + a \cdot n/n_i}.\tag{11}$$

2. Node-information only regime ($\mu = \Omega(n \cdot n_i)$). Perfect recovery is achievable if

$$\rho \leq \gamma \, n_i. \tag{12}$$

3. Synergistic regime (0 < μ < $n \cdot n_i$). Perfect recovery is achievable if

$$\rho \leq \left(\frac{1 - (a+2) p_{\max}}{\mu \ell} + \gamma\right) n_i \quad and \quad p_{\max} \leq \frac{1 - \mu \ell (\rho/n_i - \gamma)}{2 + a}. \tag{13}$$

Proof. The proof is given in Appendix B.

Theorem 4.5 unifies the recovery conditions across three regimes. In the graph-structure-only regime, the bound $p_{\max} \leq 1/(2+a\cdot n/n_i)$ shows that recovery becomes increasingly restrictive as the number of classes $K \approx n/n_i$ grows, since more communities require stronger separation. In the node-information-only regime, recovery is governed by the quality of node features, with the admissible noise bounded by $\rho \leq \gamma n_i$, scaling linearly with class size.

The synergistic regime is the most interesting: here both graph structure and node features interact, and a clear trade-off emerges between tolerance to feature noise and tolerance to misconnected edges. As seen in (13), decreasing μ from infinity (the node-information-only setting) introduces an additional term, which enlarges the admissible range of ρ and thus allows for higher noise tolerance. At the same time, from (13), p_{\max} no longer depends on the number of classes K and can be relaxed by either choosing γ close to ρ/n_i or by reducing μ . However, this improvement in p_{\max} comes at the expense of reducing the admissible range for ρ . In other words, there is a tension: one can either tolerate noisier features at the cost of stricter structural requirements, or tolerate higher misconnection rates in the graph at the cost of requiring cleaner features. This trade-off highlights the importance of balancing graph and feature contributions when designing algorithms, as exploiting both sources simultaneously provides the broadest recovery guarantees.

Algorithm 1 CADO Algorithm

324

325

326

328

330

331

332

333

334

335

336

343 344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360 361

362

364

365 366

367 368 369

370

371

372

373 374

375

376

377

- 1: Input: Number of atoms r, graph \bar{A} , node data $\{z_v\}$, step size sequence $\{\gamma_t = 2/t + 2\}$
- 2: Initialize $\{\bar{\boldsymbol{e}}_i^{(0)} \in \mathcal{U}\}_{i=1}^r, \{\bar{\boldsymbol{\theta}}_i^{(0)} \in \Theta\}_{i=1}^r$ 3: **for** $t=0,1,2,\ldots$ until convergence **do** 327
 - - 4: // Embedding Update via Conditional Gradient
 - Compute gradients $\nabla_{\bar{e}_i} \phi$ using current $\bar{\theta}_i^{(t)}$ 5:
 - Solve LMO (15); let $\tilde{E}^{(t)} = \{\tilde{e}_i^{(t)}\}$ be the solution Update: $\bar{e}_i^{(t+1)} = (1 \gamma_t)\bar{e}_i^{(t)} + \gamma_t\tilde{e}_i^{(t)}$ 6:
 - 7:
 - // Model Update via Conditional Gradient 8:
 - Compute gradients $abla_{ar{oldsymbol{ heta}}_i}\phi$ using updated $ar{oldsymbol{e}}_i^{(t+1)}$ 9:
 - Solve LMO (14); let $\tilde{\theta}_{i}^{(t)}$ be the solution 10:
 - 11: Update: $\bar{\boldsymbol{\theta}}_i^{(t+1)} = (1 \gamma_t)\bar{\boldsymbol{\theta}}_i^{(t)} + \gamma_t\tilde{\boldsymbol{\theta}}_i^{(t)}$ 12: **Return:** $\{\bar{\boldsymbol{e}}_i^{(T)}\}, \{\bar{\boldsymbol{\theta}}_i^{(T)}\}$

COMPLEXITY ANALYSIS AND ALGORITHMIC SOLUTIONS

In this section, we analyze the computational aspects of solving (6). In Appendix C, we show that an iterative procedure exists that achieves polynomial-time convergence to an ϵ -optimal solution. While theoretically appealing, this algorithm is still computationally expensive and does not scale to large graphs. To address this limitation, we propose a practical and scalable method, referred to as constrained atomic decomposition optimization solver (CADO), which solves the non-convex formulation in (4) under a fixed number of atoms. This fixed-rank approximation is well motivated by the SON regularization, which naturally promotes sparsity in the active atoms, and in practice we find that CADO consistently converges to the optimal solution.

CADO: An efficient algorithmic solution. To solve (6) with a fixed number of atoms, we propose an alternating conditional gradient algorithm, termed CADO. The algorithm alternates between updating the embedding vectors $\{\bar{e}_i\}$, which define the low-rank matrix L, and updating the cluster-level models $\{\theta_i\}$, which induce the node-specific models $\{\theta_v\}$. This alternating structure allows CADO to minimize the joint objective efficiently without explicitly enumerating the infinite set of atoms.

Each update step follows a conditional gradient (Frank-Wolfe) approach. At iteration t, a linear minimization oracle (LMO) is solved for both the embeddings and the model parameters, yielding descent directions:

$$\tilde{\boldsymbol{\theta}}_{i} = \arg \min_{\bar{\boldsymbol{\theta}}_{i} \in \Theta} \quad \left\langle \nabla_{\bar{\boldsymbol{\theta}}_{i}} \phi, \, \bar{\boldsymbol{\theta}}_{i} \right\rangle, \tag{14}$$

$$\tilde{\boldsymbol{\theta}}_{i} = \arg \min_{\bar{\boldsymbol{\theta}}_{i} \in \Theta} \quad \left\langle \nabla_{\bar{\boldsymbol{\theta}}_{i}} \phi, \, \bar{\boldsymbol{\theta}}_{i} \right\rangle, \tag{14}$$

$$\tilde{\boldsymbol{E}} = \arg \min_{\left\{\bar{\boldsymbol{e}}_{i}\right\}_{i=1}^{r}} \quad \sum_{i=1}^{r} \left\langle \nabla_{\bar{\boldsymbol{e}}_{i}} \phi, \, \bar{\boldsymbol{e}}_{i} \right\rangle, \tag{15}$$

$$\text{s.t.} \quad \boldsymbol{L} = \sum_{i=1}^r \bar{\boldsymbol{e}}_i \bar{\boldsymbol{e}}_i^\top \in \mathcal{B}, \quad \boldsymbol{\theta}_v = \sum_{i=1}^r \bar{\boldsymbol{\epsilon}}_{i,v}^2 \bar{\boldsymbol{\theta}}_i \in \Theta,$$

where $\bar{e}_{i,v}$ denotes the v-th entry of \bar{e}_i . The exact solutions of these LMOs depend on the loss function f_v and the domain Θ . In the case study presented in Section 6.1, we show that both LMOs admit efficient closed-form solutions. Detailed derivations and alrorithmic steps are provided in Appendix D.

After obtaining the LMOs, the models θ_v and embeddings \bar{e}_i are updated via a convex combination with the results from previous iterations, ensuring feasibility throughout iterations. While the alternating conditional gradient approach leads to a non-convex problem, in practice CADO converges quickly to stable solutions, as supported by our experimental results. The procedure is summarized in Algorithm 1.

6 EXPERIMENTS

For our experiments, we focus on a particular case study (Section 6.1) that admits closed-form solutions for (14) and (15). The corresponding specialized version of CADO is summarized in Algorithm 2 in Appendix D. We then evaluate our framework and CADO on node classification within this case study. Specifically, we first demonstrate the bidirectional interaction between graph structure and node information, and then empirically assess the effectiveness of CADO. Additional experiments are reported in Appendix E, further supporting our theoretical and algorithmic findings.

6.1 Case Study

Our general framework in (6) accommodates a wide range of settings through different choices of the loss functions f_v . To make the discussion concrete, we focus here on a specific case. We assume that the vectors \mathbf{z}_v are statistically independent and, within each class, identically distributed. Furthermore, in the training set, conditional on the true class label y_v , the features and the labels are independent, i.e., $p(\mathbf{x}_v, \tilde{y}_v \mid y_v = i) = p(\mathbf{x}_v \mid \mathbf{R}_i)p(\tilde{y}_v \mid \mathbf{\pi}_i)$.

We model the feature vectors as m-dimensional centered Gaussians, $\boldsymbol{x}_v \sim \mathcal{N}(0, \bar{\boldsymbol{R}}_i)$, where $\bar{\boldsymbol{R}}_i$ is the class-dependent covariance matrix. Features in each class are assumed to concentrate near a distinct linear subspace. Specifically, the eigenvalues of $\bar{\boldsymbol{R}}_i$ are divided into two groups: those larger than a fixed positive threshold ρ_+ , corresponding to a *signal subspace*, and those smaller than another threshold ρ_- , corresponding to a *noise subspace*. Naturally, we require $\rho_+ > \rho_-$.

In the training set of class i, we assume that the noisy label $\tilde{y}_v = j$ is observed with probability $\bar{\pi}_{ji}$, and define $\bar{\pi}_i = (\bar{\pi}_{ji})_j$. We further assume that $\bar{\pi}_{ii}$ is significantly larger than $\bar{\pi}_{ji}$ for $j \neq i$, i.e., the probability of observing the correct label is significantly higher than that of any incorrect label.

To instantiate (6) in this setting, we choose Θ and f_v as follows. Each model parameter θ consists of a covariance R and a label distribution π . The node-level losses are defined by

$$f_{\text{feature}}(\boldsymbol{x}; \boldsymbol{R}) \coloneqq \frac{1}{m} \left(\boldsymbol{x}^{\top} \boldsymbol{R}^{-1} \boldsymbol{x} + \text{Tr}(\boldsymbol{R}) \right),$$
 (16)

$$f_{\text{label}}(\tilde{y}=j;\boldsymbol{\pi}) := -\pi_{j}.$$
 (17)

and combined as

$$f_{v}(\boldsymbol{z};\boldsymbol{\theta}) = \begin{cases} f_{\text{feature}}(\boldsymbol{x};\boldsymbol{R}) + \beta f_{\text{label}}(\tilde{\boldsymbol{y}};\boldsymbol{\pi}) & v \in \mathcal{V}_{\text{train}} \\ f_{\text{feature}}(\boldsymbol{x};\boldsymbol{R}) & v \in \mathcal{V}_{\text{test}} \end{cases}$$
(18)

We constrain R to lie in the set $\mathbb{S}_{\rho_-,\rho_+}$ of symmetric matrices with eigenvalues in $[\rho_-,\rho_+]$, i.e.,

$$\mathbb{S}_{\rho_{-},\rho_{+}} = \{ \mathbf{R} \mid \mathbf{R} = \mathbf{R}^{\top}, \ \forall \mathbf{x} \in \mathbb{R}^{n} : \ \rho_{-} \|\mathbf{x}\|_{2}^{2} \le \mathbf{x}^{\top} \mathbf{R} \mathbf{x} \le \rho_{+} \|\mathbf{x}\|_{2}^{2} \}.$$
 (19)

Similarly, we constrain π to lie in the standard simplex $\Delta = \{\pi = (\pi_k) \mid \pi_k \geq 0, \ \sum_k \pi_k = 1\}$, Accordingly, we take $\Theta = \mathbb{S}_{\rho_-,\rho_+} \times \Delta$.

6.2 EXPERIMENTAL SETUP

Data generation model. We generate a synthetic dataset with K clusters of equal size n_0 , so that the total number of nodes is $n = Kn_0$. The graph is drawn from an SBM, while node features follow a Gaussian mixture model.

Graph structure: Edges are placed independently, with nodes in the same cluster connected with probability p and nodes in different clusters connected with probability q.

Node features: Each cluster generates zero-mean Gaussian feature vectors with a covariance matrix whose eigenvalues encode signal and noise. Specifically, m_{ω} eigenvalues are set to ω^2 (noise), while the remaining eigenvalues are set to σ^2 (signal). The parameters ω and m_{ω} jointly control the noise level and hence the degree of separability between the feature subspaces of different clusters.

Node labels: From each cluster, $n_{\rm t}$ nodes are selected as training examples. A training node is assigned the correct label with probability π and an incorrect label chosen uniformly at random otherwise. This design allows us to explore the effect of the fraction of labeled data and the level of label noise.

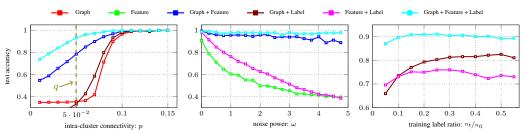


Figure 2: Left: Test accuracy vs. p; Middle: Test accuracy vs. ω ; Right: Test accuracy vs. training label ratio. All plots highlight the synergy of combining graph, features, and labels, with the highest accuracy achieved by the full framework (Graph + Feature + Label).

Hyperparameters. A full list of hyperparameters is provided in Appendix E. Unless otherwise specified, we adopt the following defaults: K=3 clusters with $n_0=300$ nodes each, and we fix the number of atoms to r=K. Edge probabilities are p=0.1 within clusters and q=0.05 across clusters. Node features are 6-dimensional, with $m_\omega=4$ noisy eigenvalues of magnitude $\omega=0.04$ in the covariance matrix, yielding well-separated feature subspaces. For training labels, we set the label ratio to $\frac{n_t}{n_0}=0.2$ and assume no label noise $(\pi=1)$. To simplify the experimental setup, we reparameterize the original scaling setup by introducing effective weights for each term. While our theoretical framework uses a scaling of 1 for the graph term, μ for the feature term, and $\mu\beta$ for the label term, we fix the weights to $\beta_g=1.0$, $\beta_f=2.5$, and $\beta_l=13.0$, corresponding to the graph, feature, and label terms, respectively.

Evaluation settings. To assess the contribution of each information source (graph structure, node features, and node labels) we perform an ablation study by selectively enabling or disabling each term in the objective. Specifically, we evaluate the graph-only setting, the feature-only setting, all pairwise combinations (graph+features, graph+levels, and features+labels), and the full model integrating all three. These ablations are implemented by setting the corresponding regularization weights to zero and solving the resulting problem with CADO, except for the graph-only case, where we apply spectral clustering instead of the ablated CADO variant.

6.3 Our Results

We present three sets of experiments in Figure 2, each isolating the impact of one component while keeping the others fixed. For each configuration, we report accuracy on test nodes.

Graph structure (left figure). We vary the intra-cluster edge probability p while keeping features and labels fixed. Our framework consistently outperforms graph-only and pairwise baselines, especially as the graph becomes less informative (i.e., $p \approx q$).

Feature quality (center figure). We degrade the signal-to-noise ratio by increasing the noise level ω in the feature covariance. Even when features become unreliable, our method remains robust by leveraging the graph structure.

Training label availability (right figure). We vary the ratio of labeled training nodes. Our framework effectively exploits even a small number of noisy labels when integrating graph and feature information.

Detailed discussions of these experiments together with additional experiments are provided in Appendix E, providing further support for our theoretical and algorithmic contributions. Overall, our results confirm that the proposed framework effectively integrates graph, feature, and label information. CADO consistently recovers the true structure when the theoretical conditions are met and remains robust as the problem becomes harder $(p\downarrow,\omega\uparrow)$. The first plot illustrates how features and labels enhance spectral clustering when the graph is weak, while the latter two show how graph structure strengthens classification when node-specific information is degraded.

Conclusion. We proposed a novel optimization framework for transductive node classification, formulated through the lens of convex graph clustering. Our approach integrates graph structure, node features, and labels within a unified convex formulation, and we developed an efficient algorithmic solution based on a fixed number of atoms. We established theoretical guarantees for perfect recovery, showing that combining graph and feature information requires strictly milder conditions than using either source alone. Experimental results corroborate our theory, demonstrating that node features enhance graph clustering while graph structure strengthens classification.

REFERENCES

- Emmanuel Abbe. Community detection and stochastic block models: recent developments. *Journal of Machine Learning Research*, 18(177):1–86, 2018.
- Emmanuel Abbe, Afonso S Bandeira, and Georgina Hall. Exact recovery in the stochastic block model. *IEEE Transactions on information theory*, 62(1):471–487, 2015.
 - Daniel Aloise, Amit Deshpande, Pierre Hansen, and Preyas Popat. Np-hardness of euclidean sum-of-squares clustering. *Machine learning*, 75(2):245–248, 2009.
 - Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. *Advances in neural information processing systems*, 14, 2001.
 - Smriti Bhagat, Graham Cormode, and S Muthukrishnan. Node classification in social networks. *Social network data analytics*, pp. 115–148, 2011.
 - Badri Narayan Bhaskar, Gongguo Tang, and Benjamin Recht. Atomic norm denoising with applications to line spectral estimation. *IEEE Transactions on Signal Processing*, 61(23):5987–5999, 2013.
 - Norbert Binkiewicz, Joshua T Vogelstein, and Karl Rohe. Covariate-assisted spectral clustering. *Biometrika*, 104(2):361–377, 2017.
 - Guillaume Braun and Masashi Sugiyama. Vec-sbm: Optimal community detection with vectorial edges covariates. In *International Conference on Artificial Intelligence and Statistics*, pp. 532–540. PMLR, 2024.
 - Guillaume Braun, Hemant Tyagi, and Christophe Biernacki. An iterative clustering algorithm for the contextual stochastic block model with optimality guarantees. In *International Conference on Machine Learning*, pp. 2257–2291. PMLR, 2022.
 - Yudong Chen, Ali Jalali, Sujay Sanghavi, and Huan Xu. Clustering partially observed graphs via convex optimization. *The Journal of Machine Learning Research*, 15(1):2213–2238, 2014.
 - Eric C Chi and Stefan Steinerberger. Recovering trees with convex clustering. *SIAM Journal on Mathematics of Data Science*, 1(3):383–407, 2019.
 - Julien Chiquet, Pierre Gutierrez, and Guillem Rigaill. Fast tree inference with weighted fusion penalties. *Journal of Computational and Graphical Statistics*, 26(1):205–216, 2017.
 - Petr Chunaev. Community detection in node-attributed social networks: a survey. *Computer Science Review*, 37:100286, 2020.
 - Yash Deshpande, Subhabrata Sen, Andrea Montanari, and Elchanan Mossel. Contextual stochastic block models. *Advances in Neural Information Processing Systems*, 31, 2018.
 - Maximilien Dreveton, Felipe Fernandes, and Daniel Figueiredo. Exact recovery and bregman hard clustering of node-attributed stochastic block model. *Advances in Neural Information Processing Systems*, 36:37827–37848, 2023.
 - Bruce Hajek, Yihong Wu, and Jiaming Xu. Achieving exact cluster recovery threshold via semidefinite programming. *IEEE Transactions on Information Theory*, 62(5):2788–2797, 2016.
 - Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
 - Toby Dylan Hocking, Armand Joulin, Francis Bach, and Jean-Philippe Vert. Clusterpath an algorithm for clustering using convex fusion penalties. In 28th international conference on machine learning, pp. 1, 2011.
 - Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.

- Yaofang Hu and Wanjie Wang. Network-adjusted covariates for community detection. *Biometrika*, 111(4):1221–1240, 2024.
- Adam Jaeger and David Banks. Cluster analysis: A modern statistical review. *Wiley Interdisciplinary Reviews: Computational Statistics*, 15(3):e1597, 2023.
 - Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
 - Ramya Korlakai Vinayak, Samet Oymak, and Babak Hassibi. Graph clustering with missing data: Convex algorithms and analysis. *Advances in Neural Information Processing Systems*, 27, 2014.
 - Xiaodong Li, Yudong Chen, and Jiaming Xu. Convex relaxation methods for community detection. *Statistical science*, 36(1):2–15, 2021. ISSN 0883-4237.
 - Fredrik Lindsten, Henrik Ohlsson, and Lennart Ljung. Clustering using sum-of-norms regularization: With application to particle filter output computation. In 2011 IEEE Statistical Signal Processing Workshop (SSP), pp. 201–204. IEEE, 2011.
 - Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*, pp. 43–52, 2015.
 - Miller McPherson, Lynn Smith-Lovin, and James M Cook. Birds of a feather: Homophily in social networks. *Annual review of sociology*, 27(1):415–444, 2001.
 - Galileo Namata, Ben London, Lise Getoor, Bert Huang, and U Edu. Query-driven active surveying for collective classification. In *International workshop on mining and learning with graphs (MLG)*, volume 8, 2012.
 - Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 14, 2001.
 - Ashkan Panahi, Devdatt Dubhashi, Fredrik D Johansson, and Chiranjib Bhattacharyya. Clustering by sum of norms: Stochastic incremental algorithm, convergence and cluster recovery. In *International conference on machine learning*, pp. 2769–2777. PMLR, 2017.
 - Nikhil Rao, Parikshit Shah, and Stephen Wright. Forward–backward greedy algorithms for atomic norm regularization. *IEEE Transactions on Signal Processing*, 63(21):5798–5811, 2015.
 - Satu Elisa Schaeffer. Graph clustering. Computer science review, 1(1):27–64, 2007.
 - Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3), 2008.
 - Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
 - Mahdi Soltanolkotabi and Emmanuel J Candes. A geometric analysis of subspace clustering with outliers. *Annals of Statistics*, 2012.
 - Defeng Sun, Kim-Chuan Toh, and Yancheng Yuan. Convex clustering: Model, theoretical guarantee and efficient algorithm. *Journal of Machine Learning Research*, 22(9):1–32, 2021.
 - Kean Ming Tan and Daniela Witten. Statistical properties of convex clustering. *Electronic journal of statistics*, 9(2):2324, 2015.
 - Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018.
 - Yihong Wu, Jiaming Xu, and Bruce Hajek. Achieving exact cluster recovery threshold via semidefinite programming under the stochastic block model. In 2015 49th Asilomar Conference on Signals, Systems and Computers, pp. 1070–1074. IEEE, 2015.

Rui Xu and Donald Wunsch. Survey of clustering algorithms. IEEE Transactions on neural networks, 16(3):645–678, 2005. Bowei Yan and Purnamrita Sarkar. Covariate regularized community detection in sparse graphs. Journal of the American Statistical Association, 116(534):734–745, 2021. Shenghao Yang and Kimon Fountoulakis. Weighted flow diffusion for local graph clustering with node attributes: An algorithm and statistical guarantees. In International Conference on Machine Learning, pp. 39252–39276. PMLR, 2023. Changbo Zhu, Huan Xu, Chenlei Leng, and Shuicheng Yan. Convex optimization procedure for clustering: Theoretical revisit. Advances in Neural Information Processing Systems, 27, 2014. Xiaojin Zhu, Zoubin Ghahramani, and John Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. ICML, 3:912-919, 2003.

MATHEMATICAL NOTATION.

 To simplify notation and enhance readability, we use the following conventions throughout the paper. Collections such as $\{\lambda_a\}_{a\in\mathcal{A}}$, $\{L_{vv}\}_{v\in\mathcal{V}}$, $\{L_{uv}\}_{u,v\in\mathcal{V}}$, and $\{\theta_v\}_{v\in\mathcal{V}}$ are often abbreviated to $\{\lambda_a\}$, $\{L_{vv}\}$, $\{L_{uv}\}$, and $\{\theta_v\}$ when the indexing set is clear from context. When two nodes u and v belong to the same cluster, we write $u,v\in\mathcal{C}_i$; otherwise, we denote $u\in\mathcal{C}_i$ and $v\in\mathcal{C}_j$ to indicate distinct clusters.

We denote vectors by bold lowercase letters (e.g., \boldsymbol{x}) and matrices by bold uppercase letters (e.g., \boldsymbol{W}). The entry in the v^{th} row and u^{th} column of a matrix \boldsymbol{W} is denoted by W_{vu} , and its transpose by \boldsymbol{W}^{\top} . The Euclidean inner product between vectors is written as $\langle \cdot, \cdot \rangle$, and the matrix inner product is defined by $\langle \boldsymbol{A}, \boldsymbol{C} \rangle := \operatorname{tr}(\boldsymbol{A}\boldsymbol{C}^{\top})$, where $\operatorname{tr}(\cdot)$ denotes the trace operator.

The indicator function $I_{\Theta}(\theta)$ refers to the indicator of the feasible set Θ , and $\partial I_{\Theta}(\theta)$ denotes the normal cone of Θ at point θ . The cardinality of a set T is denoted by |T|. We use $\mathbf{1}_n$ and $\mathbf{0}_n$ to denote the n-dimensional all-one and all-zero vectors, respectively. An $m \times n$ all-zero matrix is denoted by $O_{m \times n}$, with subscripts omitted when dimensions are clear from context.

A ATOMIC NORM REVIEW

 In many signal processing and machine learning applications, the objective is to reconstruct a signal that admits a simple or structured representation. The concept of *atomic norms* provides a unifying framework for this purpose by promoting simplicity in the signal representation. Specifically, an atomic norm can be used as a convex regularizer that induces a desired structure, such as sparsity or low-rankness, by relying on a predefined set of fundamental building blocks called *atoms*.

A.1 ATOMS AND ATOMIC NORMS

The basic elements used to represent a signal are referred to as *atoms*. For a signal x, the atomic set \mathcal{A} is a collection of these basic elements, allowing the signal to be expressed as a nonnegative linear combination of a small number of atoms from \mathcal{A} . The atomic norm, denoted by $||x||_{\mathcal{A}}$, quantifies the complexity of the signal in terms of how economically it can be represented using these atoms.

Formally, the atomic norm is defined as:

$$\|\boldsymbol{x}\|_{\mathcal{A}} = \inf \left\{ \sum_{i} c_{i} : \boldsymbol{x} = \sum_{i} c_{i} \boldsymbol{a}_{i}, \ \boldsymbol{a}_{i} \in \mathcal{A}, \ c_{i} \geq 0 \right\}.$$
 (20)

This norm is often employed as a regularizer in optimization problems to encourage solutions that are structurally simple, such as sparse vectors or low-rank matrices. Different choices of the atomic set \mathcal{A} yield different atomic norms. Two important examples are as follows:

Sparsity (\ell_1-norm): For sparse signals, the atomic set consists of the signed canonical basis vectors:

$$\mathcal{A}_{\ell_1} = \{ \pm \boldsymbol{e}_i \mid i = 1, \dots, n \},\,$$

where e_i denotes the *i*-th canonical basis vector in \mathbb{R}^n . The induced atomic norm is:

$$\|\boldsymbol{x}\|_{\mathcal{A}_{\ell_1}} = \sum_{i=1}^n |x_i| = \|\boldsymbol{x}\|_1,$$

which is the familiar ℓ_1 -norm, widely used to promote sparsity.

• Low-rank matrices (nuclear norm): For low-rank matrix recovery, such as in matrix completion or graph clustering, the atomic set consists of rank-one matrices with unit norm:

$$\mathcal{A}_* = \{ \mathbf{u}\mathbf{v}^\top \mid \mathbf{u} \in \mathbb{R}^m, \, \mathbf{v} \in \mathbb{R}^n, \, \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1 \}.$$

The induced atomic norm is:

$$\|\mathbf{X}\|_{\mathcal{A}_*} = \sum_{i=1}^{\min(m,n)} \sigma_i(\mathbf{X}) = \|\mathbf{X}\|_*,$$

which is the nuclear norm, equal to the sum of the singular values of X. This norm is the convex surrogate of the rank function and promotes low-rank structure.

For symmetric matrices, the above definitions can be simplified. If $X = X^{T}$, we can define the symmetric atomic sets

$$\mathcal{A}_{\mathrm{sym},+} \ = \ \{ \, \mathbf{u} \mathbf{u}^\top : \| \mathbf{u} \|_2 = 1 \, \}, \qquad \mathcal{A}_{\mathrm{sym},\pm} \ = \ \{ \, \pm \mathbf{u} \mathbf{u}^\top : \| \mathbf{u} \|_2 = 1 \, \}.$$

Then, with the eigenvalue decomposition $\mathbf{X} = \sum_i \lambda_i \, \mathbf{u}_i \mathbf{u}_i^{\top}$ (like graph Laplacian matrix):

$$\|\mathbf{X}\|_{\mathcal{A}_{\mathrm{sym},+}} = \sum_{i} \lambda_{i} = \mathrm{tr}(\mathbf{X}) \text{ for } \mathbf{X} \succeq 0,$$

and, for general symmetric X,

$$\|\mathbf{X}\|_{\mathcal{A}_{\mathrm{sym},\pm}} \ = \ \sum_i |\lambda_i| \ = \ \|\mathbf{X}\|_*.$$

In particular, for symmetric matrices the nuclear norm equals the sum of absolute eigenvalues, and for symmetric PSD matrices it reduces to the trace.

B PROOF OF THEOREM 4.5

We begin with an overview of the proof methodology, followed by the problem setup and optimization formulation. We then outline the key elements, focusing on optimality conditions for exact recovery. The main proof is divided into three parts, each covering a different information regime: (i) node-specific, (ii) graph-based, and (iii) combined. For each, we construct a dual certificate via a "guess-and-golfing" approach and establish recovery conditions. We conclude by stating the theorem.

B.1 Proof Overview

Our proof follows a standard dual certificate approach. We aim to verify that the ideal solution—comprising one atom per cluster—satisfies the optimality (KKT) conditions of the convex problem (6). The core idea is to construct a dual certificate matrix \mathbf{Z} (and associated subgradients $\{\mathbf{g}_{vu}\}$) such that the ideal solution is locally optimal.

To achieve this, we employ a *guess-and-golfing strategy*. Starting from a natural or trivial guess for \mathbf{Z} (e.g., diagonal or block-structured), we iteratively refine it to meet the necessary conditions. Each refinement step aims to adjust the certificate to satisfy one or more of the following:

- sign constraints ensuring dual feasibility;
- inner product constraints capturing primal-dual consistency;
- and positive semidefiniteness of the gap matrix enforcing inequality optimality.

B.2 PROBLEM SETUP

Optimization Problem Restatement. The optimization problem we analyze, with the goal of establishing conditions for perfect recovery, is our convex framework augmented with a sum-of-norms (SON) regularization term. For clarity and completeness, we restate it below, as in equation (6).

$$\min_{\mathbf{U} \in \mathcal{U}, \{\lambda_{a} \ge 0\}_{a \in \mathcal{A}}} \quad \phi(\mathbf{U}) + \mu_{0} \sum_{\mathbf{a} \in \mathcal{A}} \lambda_{\mathbf{a}} + \mu_{1} R(\mathbf{U})$$
s.t.
$$\mathbf{U} = \sum_{\mathbf{a} \in \mathcal{A}} \lambda_{\mathbf{a}} \mathbf{a}$$
(21)

In this formulation, $\boldsymbol{U}:=(\boldsymbol{L},\{\boldsymbol{\theta}_v\}_{v\in\mathcal{V}})$. The feasible set is $\mathcal{U}:=\mathcal{B}\times\Theta^{|\mathcal{V}|}$, where \mathcal{B} denotes the set of symmetric matrices with entries in [0,1] and unit diagonal, and Θ is the feasible domain for node-specific models. The objective $\phi(\boldsymbol{U})$ is defined as: $\phi(\boldsymbol{U}):=-\langle \bar{\boldsymbol{A}},\boldsymbol{L}\rangle+\mu\sum_{v\in\mathcal{V}}f_v(\boldsymbol{z}_v;\boldsymbol{\theta}_v),$ and $R(\boldsymbol{U})$ is the SON regularization term. Each atom $\boldsymbol{a}\in\mathcal{A}$ is of the form $\boldsymbol{a}=(\boldsymbol{e}\boldsymbol{e}^{\top},\{\epsilon_v^2\boldsymbol{\theta}\}_{v\in\mathcal{V}}),$ and comes from the atomic dictionary \mathcal{A} defined as:

$$\mathcal{A} := \left\{ \left(e e^{\top}, \{ \epsilon_v^2 \boldsymbol{\theta} \}_{v \in \mathcal{V}} \right) \mid \boldsymbol{e} = (\epsilon_v)_{v \in \mathcal{V}}, \ \| \boldsymbol{e} \| = 1, \ \boldsymbol{\theta} \in \Theta \right\}. \tag{22}$$

Using the definitions above, the optimization problem (21) can be explicitly given as:

$$\min_{\substack{\boldsymbol{L} \in \mathcal{B}, \ \{\boldsymbol{\theta}_{v} \in \Theta\}_{v \in \mathcal{V}}, \\ \{\lambda_{\boldsymbol{a}} \geq 0\}_{\boldsymbol{a} \in \mathcal{A}}}} -\langle \bar{\boldsymbol{A}}, \boldsymbol{L} \rangle + \mu \sum_{v} f_{v}(\boldsymbol{z}_{v}, \boldsymbol{\theta}_{v}) + \mu_{0} \sum_{\boldsymbol{a} \in \mathcal{A}} \lambda_{\boldsymbol{a}} + \mu_{1} \sum_{u < v} \|\boldsymbol{\theta}_{v} - \boldsymbol{\theta}_{u}\|$$
s.t.
$$\boldsymbol{L} = \sum_{a \in \mathcal{A}} \lambda_{a} \boldsymbol{e}_{a} \boldsymbol{e}_{a}^{\top}, \qquad \boldsymbol{\theta}_{v} = \sum_{a \in \mathcal{A}} \lambda_{a} \epsilon_{a,v}^{2} \boldsymbol{\theta}_{a}$$

$$(23)$$

Equivalent Optimization Problem. We simplify the problem in (23) by enforcing $L_{vv}=1$ for all v, which implies that the sum of coefficients λ_a is fixed, and the conditions $L_{uv} \leq 1$ and $\theta_v \in \Theta$ are automatically satisfied (see Lemma B.2). We also eliminate the auxiliary variables θ_v by expanding their definition in the objective, leading to the following simplified problem:

$$\min_{\substack{\{L_{uv} \geq 0\}_{u,v,} \\ \{\lambda_a \geq 0\}_{a \in \mathcal{A}}}} -\langle \bar{\boldsymbol{A}}, \boldsymbol{L} \rangle + \mu \sum_{v \in \mathcal{V}} f_v \left(\boldsymbol{z}_v, \sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,v}^2 \boldsymbol{\theta}_a \right) + \mu_1 \sum_{u < v} \left\| \sum_{a \in \mathcal{A}} \lambda_a \left(\epsilon_{a,u}^2 - \epsilon_{a,v}^2 \right) \boldsymbol{\theta}_a \right\|$$
s.t.
$$\boldsymbol{L} = \sum_{a \in \mathcal{A}} \lambda_a \boldsymbol{e}_a \boldsymbol{e}_a^\top, \qquad L_{vv} = 1, \quad \forall v \in \mathcal{V}.$$
(24)

B.3 Proof Elements

Ideal Solution. We define the ideal (perfect recovery) solution as one consisting of K atoms—one per cluster—denoted by $\{(\lambda_i^*, e_i^*, \theta_i^*)\}_{i=1}^K$. These ideal atoms take the form:

$$\lambda_i^* = n_i, \qquad \epsilon_{i,v}^* = \begin{cases} \frac{1}{\sqrt{n_i}}, & v \in \mathcal{C}_i \\ 0, & \text{otherwise} \end{cases}, \qquad \boldsymbol{e}_i^* = (\epsilon_{i,v}^*)_{v \in \mathcal{V}}. \tag{25}$$

$$\{\boldsymbol{\theta}_i^*\}_{i=1}^K = \arg\min_{\{\boldsymbol{\theta}_i \in \Theta\}_i} \sum_{v \in \mathcal{V}} h_v(\boldsymbol{\theta}_{y_v})$$
 (26)

where $h_v(\boldsymbol{\theta}_i) = f_v(\boldsymbol{z}_v; \boldsymbol{\theta}_i) + \gamma \sum_{i < j} n_j \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j^*\|$. This yields an ideally partitioned matrix \boldsymbol{L} with:

$$L_{uv}^* = \begin{cases} 1, & u, v \in \mathcal{C}_i \text{ for some } i \\ 0, & \text{otherwise} \end{cases}$$
 (27)

and the following characteristic condition derives from the first-order optimality condition of (26):

$$\frac{1}{n_i} \sum_{v \in \mathcal{C}_i} \nabla h_v(\boldsymbol{\theta}_{y_v}^*) \coloneqq \boldsymbol{\omega}_i \in -\partial I_{\Theta}(\boldsymbol{\theta}_i^*)$$
 (28)

where $\partial I_{\Theta}(\theta_i^*)$ is the normal cone of Θ at θ_i^* defined as

$$\partial I_{\Theta}(\boldsymbol{\theta}_{i}^{*}) := \begin{cases} \{\boldsymbol{\omega} \in \mathbb{R}^{d} : \langle \boldsymbol{\omega}, \boldsymbol{\theta} - \boldsymbol{\theta}_{i}^{*} \rangle \leq 0, \ \forall \boldsymbol{\theta} \in \Theta \} & \text{if} \quad \boldsymbol{\theta}_{i}^{*} \in \Theta \\ \varnothing & \text{if} \quad \boldsymbol{\theta}_{i}^{*} \notin \Theta \end{cases}$$
(29)

Note that,

$$\nabla h_v(\boldsymbol{\theta}_i^*) = \nabla f_v(\boldsymbol{\theta}_i^*) + \gamma \sum_{u} \mathbf{g}_{vu}.$$
(30)

with $\gamma := \mu_1/\mu$. The terms \mathbf{g}_{vu} in (30) correspond to subgradients of the SON regularization term, and can be calculated as:

$$\mathbf{g}_{vu} = \begin{cases} \underset{u}{\text{any }} \mathbf{g}_{vu} \text{ with } \|\mathbf{g}_{vu}\| \leq 1, & \text{if } u, v \in \mathcal{C}_i \\ \frac{\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*}{\|\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*\|}, & \text{if } u \in \mathcal{C}_i, \ v \in \mathcal{C}_j, \ i \neq j. \end{cases}$$
(31)

By symmetry of the SON term, it is required that $\mathbf{g}_{vu} = -\mathbf{g}_{uv}$. As a results, its easy to show that

$$\omega_i = \nabla F_i(\boldsymbol{\theta}_i^*) + \gamma \sum_{i < j} n_j \frac{\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*}{\|\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*\|}$$
(32)

where $\nabla F_i(\boldsymbol{\theta}) := \frac{1}{n_i} \sum_{v \in \mathcal{C}_i} \nabla f_v(\boldsymbol{z}_v; \boldsymbol{\theta}).$

Optimality Conditions. To establish that the ideal solution is optimal for the simplified problem (24), it suffices to verify the first-order optimality (KKT) conditions *at the ideal solution*. Specifically, we aim to construct a dual certificate: a symmetric matrix $\mathbf{Z} \in \mathbb{R}^{n \times n}$ and a collection of subgradients $\{\mathbf{g}_{vu}\}$ such that the following optimality conditions are satisfied:

$$(\mathbf{Z} - \bar{\mathbf{A}})_{uv} \le 0,$$
 if $u, v \in \mathcal{C}_i,$ (33)

$$(\mathbf{Z} - \bar{\mathbf{A}})_{uv} \ge 0,$$
 if $u \in \mathcal{C}_i, v \in \mathcal{C}_j, i \ne j,$ (34)

$$\mathbf{e}_{i}^{*\top} (\mu \mathbf{D}(\mathbf{\theta}_{i}^{*}) - \mathbf{Z}) \mathbf{e}_{i}^{*} = 0, \qquad \forall i, \tag{35}$$

$$e^{\top}(\mu \mathbf{D}(\boldsymbol{\theta}) - \mathbf{Z})e > 0, \qquad \forall (e, \boldsymbol{\theta}) \in \mathcal{A}.$$
 (36)

where $D(\theta) := \operatorname{diag}(d_v(\theta))$ with

$$d_v(\boldsymbol{\theta}) = \langle \nabla h_v(\boldsymbol{\theta}_{u_v}^*), \ \boldsymbol{\theta} \rangle \ . \tag{37}$$

Conditions (33) and (34) ensures complementary slackness with respect to the matrix variable L. Also, importantly, the diagonal elements of \mathbf{Z} are unconstrained and need not satisfy the sign conditions, as they are not involved in the objective due to the fixed diagonal constraint $L_{vv} = 1$.

Inequality (36) must hold for all atoms $(e, \theta) \in A$, while equality (35) holds for the K ideal atoms $\{(e_i^*, \theta_i^*)\}_{i=1}^K$.

Hence, the KKT conditions reduce to constructing a dual matrix \mathbf{Z} and a consistent set of subgradients $\{\mathbf{g}_{vu}\}$ that jointly satisfy equations (33)–(36). Note that \mathbf{g}_{vu} is undefined when $u,v\in\mathcal{C}_i$ (i.e., both nodes belong to the same cluster). In this case, \mathbf{g}_{vu} acts as a free design parameter and must satisfy the subgradient constraint $\|\mathbf{g}_{vu}\| \leq 1$. We jointly construct the collection $\{\mathbf{g}_{vu}\}_{u,v\in\mathcal{C}_i}$ and the dual certificate \mathbf{Z} such that the full set of optimality conditions (33)–(36) are all satisfied.

B.4 PERFECT RECOVERY WITH COMBINED GRAPH AND NODE-SPECIFIC INFORMATION

The goal of this section is to demonstrate that by *combining graph structure and node-specific information*, we obtain *looser requirements* on the parameters λ , μ , and ρ for perfect recovery.

B.4.1 CONSTRUCTING Z

We begin by constructing a candidate for the dual certificate \mathbf{Z} and the subgradients $\{\mathbf{g}_{vu}\}$ that satisfy the optimality conditions (33)- (36). We proceed via a step-by-step design, where we iteratively refine a candidate \mathbf{Z} to satisfy each condition.

Step 1: Structured Guess. We begin with a matrix $S \in \mathbb{R}^{n \times n}$ defined blockwise according to:

$$S_{uv} := \begin{cases} 0 & \text{if } u,v \in \mathcal{C}_i \text{ and } \{u,v\} \in E, \\ -a_{ii} & \text{if } u,v \in \mathcal{C}_i \text{ and } \{u,v\} \notin E, \\ 0 & \text{if } u \in \mathcal{C}_i, \ v \in \mathcal{C}_j, \ \{u,v\} \notin E, \\ a_{ij} & \text{if } u \in \mathcal{C}_i, \ v \in \mathcal{C}_j, \ \{u,v\} \in E, \end{cases}$$

where a_{ii} , $a_{ij} > 0$ are scalar parameters. This form penalizes incorrect intra-cluster disconnections and inter-cluster connections.

Step 2: Projection for Orthogonality. To ensure that \mathbf{Z} is orthogonal to the ideal atoms (i.e., satisfies (35)), we project \mathbf{S} onto the orthogonal complement of the span of ideal atoms:

$$\mathbf{Z}_{s} := \mathbf{P}^{\perp}(\mathbf{S} - \lambda \mathbf{I})\mathbf{P}^{\perp}, \text{ where } \mathbf{P}^{\perp} := \mathbf{I} - \boldsymbol{E}\boldsymbol{E}^{\top},$$

and $E = [e_1^*, \dots, e_K^*]$ is the matrix of ideal cluster indicators. This construction ensures $e_i^{*\top} \mathbf{Z}_{s} e_i^* = 0$.

Step 3: Feature Guess. To account the effect of $\mu D(\theta_i^*)$ we will add the term $\mu \bar{D} := \mu \operatorname{diag} \left(d_v(\theta_{y_v}^*) \right)$ to \mathbf{Z} with $d_v(\boldsymbol{\theta})$ is given in equation (37). This ensures that $\left(D(\theta_i^*) - \bar{D} \right)_{vv} = 0 \quad \forall v, i$

The dual certificate defined as:

$$\mathbf{Z} = \mathbf{Z}_{\mathrm{s}} + \mu \bar{\mathbf{D}},\tag{38}$$

B.4.2 VERIFYING OPTIMALITY CONDITION

The optimality condition (35) holds due to D's diagonal structure and the projection \mathbf{P}^{\perp} .

B.4.3 VERIFYING SIGN CONDITIONS

Next, we enforce the sign constraints (33)-(34) by analyzing the entrywise form of $\mathbf{Z} - \bar{\mathbf{A}}$. Using Assumptions 4.1–4.2, we compute worst-case deviations after projection and derive bounds:

$$\frac{1 + \lambda/n_i}{1 - \rho_{ii} - 2\delta} \le a_{ii} \le \frac{1 - \lambda/n_i}{\rho_{ii} + 2\delta},$$
$$\frac{1}{1 - \rho_{ij} - 2\delta} \le a_{ij} \le \frac{1}{\rho_{ij} + 2\delta}.$$

To ensure that the sign conditions remain valid, we must enforce $\frac{1+\lambda/n_i}{1-\rho_{ij}-2\delta} \leq \frac{1-\lambda/n_i}{\rho_{ij}+2\delta}$ which lead to

$$\lambda \le (1 - 2p_{\text{max}}) \, n_i \tag{39}$$

where $p_{\max} := \rho^+ + \delta$ with $\rho^+ = \max_{i,j} \rho_{ij}^+$.

B.4.4 Verifying Positive Definiteness.

Finally, to ensure $\mathbf{Z} \succeq 0$ as required by (36), we must choose λ large enough so that $\mathbf{P}^{\perp}(\lambda \mathbf{I} - \mathbf{S})\mathbf{P}^{\perp}$ has all non-positive eigenvalues.

Step 1: Design g_{vu} . To enforce condition (36), first notice

$$\tilde{d}_{v}(\boldsymbol{\theta}) := \left(\boldsymbol{D}(\boldsymbol{\theta}) - \bar{\boldsymbol{D}}\right)_{vv} = \left(d_{v}(\boldsymbol{\theta}) - d_{v}(\boldsymbol{\theta}_{y_{v}}^{*})\right) \\
= \left\langle \nabla h_{v}(\boldsymbol{\theta}_{y_{v}}^{*}), \ \boldsymbol{\theta} - \boldsymbol{\theta}_{y_{v}}^{*} \right\rangle \\
= \left\langle \nabla f_{v}(\boldsymbol{\theta}_{y_{v}}^{*}) + \gamma \sum_{u} \mathbf{g}_{vu}, \ \boldsymbol{\theta} - \boldsymbol{\theta}_{y_{v}}^{*} \right\rangle \tag{40}$$

By defining \mathbf{g}_{vu} as

$$\mathbf{g}_{vu} = \begin{cases} \frac{\nabla f_u(\boldsymbol{\theta}_i^*) - \nabla f_v(\boldsymbol{\theta}_i^*)}{\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*}, & \text{if } u, v \in \mathcal{C}_i, \\ \frac{\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_j^*}{\|\boldsymbol{\theta}_i^* - \boldsymbol{\theta}_i^*\|}, & \text{otherwise.} \end{cases}$$
(41)

we have the following result

$$\tilde{d}_v(\boldsymbol{\theta}) = \left\langle \boldsymbol{\tau}_v + \boldsymbol{\omega}_{y_v}, \ \boldsymbol{\theta} - \boldsymbol{\theta}_{y_v}^* \right\rangle \ge 0.$$
 (42)

where $\tau_{v} = \left(1 - \frac{\gamma n_{i}}{\rho}\right) \left(\nabla f_{v}(\boldsymbol{\theta}_{i}) - \nabla F_{i}(\boldsymbol{\theta}_{i})\right)$

Step 2: Bounding Node-specific Term. Substituting the dual certificate **Z**, we get:

$$\boldsymbol{\epsilon}^{\top} (\boldsymbol{D}(\boldsymbol{\theta}) - \mathbf{Z}) \boldsymbol{\epsilon} = \boldsymbol{\epsilon}^{\top} \left[\mathbf{Z}_{s} + \mu \tilde{\boldsymbol{D}}(\boldsymbol{\theta}) \right] \boldsymbol{\epsilon}$$
 (43)

where $\tilde{D}(\theta) = \operatorname{diag}\left(\tilde{d}_v(\theta)\right)$ with $\tilde{d}_v(\theta) = \langle \boldsymbol{\tau}_v + \boldsymbol{\omega}_{y_v}, \theta - \boldsymbol{\theta}_{y_v}^* \rangle$. According to (28) we know that $\langle \boldsymbol{\omega}_{y_v}, \theta - \boldsymbol{\theta}_{y_v}^* \rangle \geq 0 \quad \forall \, \theta \in \Theta$. Moreover, under assumption 4.3, there exists at most one index i for which

$$\langle \boldsymbol{\omega}_{y_v}, \boldsymbol{\theta} - \boldsymbol{\theta}_{y_v}^* \rangle \leq R.$$

Without loss of generality, assume this is index i = 0. We define the diagonal matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ by:

$$\mathbf{Q} = \mathrm{diag}(r_v) \quad ext{with} \quad r_v = egin{cases} 0 & v \in \mathcal{C}_i \ R & v
otin \mathcal{C}_i \end{cases}$$

Moreover, Based on assumption 4.3, the feasible set Θ is bounded as:

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}_i\| < \ell, \quad \forall i,$$

Using the Cauchy-Schwarz inequality, we get:

$$\langle \boldsymbol{\tau}_{n}, \boldsymbol{\theta} - \boldsymbol{\theta}_{i} \rangle > -\zeta \mathbf{I}.$$
 (44)

where $\zeta := \ell \left(\rho - \gamma \cdot n_i \right)$

Thus,

$$\tilde{D}(\theta) \succeq \mathbf{Q} - \zeta \mathbf{I}.$$
 (45)

Step 3: PSD Conditions. Putting (45) into (43), we obtain:

$$\boldsymbol{\epsilon}^{\top} (\boldsymbol{D}(\boldsymbol{\theta}) - \mathbf{Z}) \boldsymbol{\epsilon} \ge \boldsymbol{\epsilon}^{\top} (\mu \mathbf{Q} - \mu \zeta \mathbf{I} - \mathbf{Z}_{s}) \boldsymbol{\epsilon}.$$
 (46)

We want the matrix

$$\mathbf{\Gamma} := \mathbf{P}^{\perp}(\lambda \mathbf{I} - \mathbf{S})\mathbf{P}^{\perp} + \mu \mathbf{Q} - \mu \zeta \mathbf{I}$$

to be positive semi-definite. By invoking Lemma B.3, we obtain sufficient conditions on λ, μ , and ρ that guarantee positive semi-definiteness, and hence, perfect recovery. To apply Lemma B.3, we define the following block components:

$$\Gamma_{11} = \left(\mathbf{P}^{\perp}(\lambda \mathbf{I} - \mathbf{S})\mathbf{P}^{\perp}\right)_{(\mathcal{C}_i, \mathcal{C}_i)} - \mu \zeta \mathbf{I}_{n_i} \tag{47}$$

$$\Gamma_{22} = \left(\mathbf{P}^{\perp}(\lambda \mathbf{I} - \mathbf{S})\mathbf{P}^{\perp}\right)_{(\mathcal{C}_{\neq i}, \mathcal{C}_{\neq i})} + \mu R \mathbf{I}_{n-n_i} - \mu \zeta \mathbf{I}_{n-n_i}$$
(48)

$$\Gamma_{12} = \left(-\mathbf{P}^{\perp}\mathbf{S}\mathbf{P}^{\perp}\right)_{(\mathcal{C}_{i},\mathcal{C}_{\neq i})} \tag{49}$$

where for a matrix \mathbf{M} , the notation $\mathbf{M}_{(\mathcal{A},\mathcal{B})}$ denotes the submatrix with rows indexed by \mathcal{A} and columns indexed by \mathcal{B} . Here, $\mathcal{C}_{\neq i} := \mathcal{V} \setminus \mathcal{C}_i$ is the set of nodes not belonging to \mathcal{C}_i .

Based on the block definitions in (49) and applying Lemma B.6, we can compute the quantities α , β , and $\sigma_{\max}^2(\Gamma_{12})$ required in Lemma B.3 as follows:

$$\alpha = \lambda - a \cdot p_{\text{max}} \cdot n_i - \mu \zeta \tag{50}$$

$$\beta = \lambda - a \cdot p_{\text{max}} \cdot (n - n_i) + \mu R - \mu \zeta \tag{51}$$

$$\sigma_{\max}^2(\mathbf{\Gamma_{12}}) = a^2 \cdot p_{\max}^2 \cdot n_i \cdot (n - n_i)$$
(52)

For positive semidefiniteness, the condition $\sigma_{\max}^2(\Gamma_{12}) \leq \alpha\beta$ must hold. Substituting the above expressions, this reduces to verifying

$$\lambda \cdot (\lambda - a \cdot p_{\max} \cdot n - \mu \zeta) + \mu \cdot R \cdot (\lambda - a \cdot p_{\max} \cdot n_i - \mu \zeta) \ge 0$$
 (53)

B.4.5 Combining Conditions

The sufficient conditions derived earlier can be unified by analyzing three distinct regimes, depending on whether the graph structure dominates, the node-specific information dominates, or both contribute jointly.

Case 1: Graph structure dominates. When $\mu \to 0$, only the graph structure contributes. In this case, the condition for positive semidefiniteness reduces to

$$\lambda \ge a \cdot p_{\max} n. \tag{54}$$

On the other hand, from the spectral bound (39), we also require

$$\lambda \le (1 - 2p_{\max}) n_i. \tag{55}$$

Combining the two yields the feasibility condition

$$p_{\max} \le \frac{1}{2 + a \cdot n/n_i}.\tag{56}$$

Thus, recovery is possible only when the graph is sufficiently assortative: dense within clusters and sparse across clusters. The threshold depends on the ratio n/n_i , which can be approximated by the number of classes K.

Case 2: Node-specific information dominates. When $\mu \gg -\langle \bar{A}, L \rangle$, in particular when $\mu = \Omega(n \cdot n_i)$, the node-specific information outweighs the graph structure. In this case we require

$$\lambda \geq a \cdot p_{\max} n_i + \mu \zeta. \tag{57}$$

On the other hand, the spectral bound (39) still enforces

$$\lambda \le (1 - 2p_{\text{max}}) n_i. \tag{58}$$

Combining the two yields the feasibility condition

$$\rho \le \left(\frac{1 - (a+2) p_{\text{max}}}{\mu \ell} + \gamma\right) n_i \approx \gamma n_i, \tag{59}$$

showing that perfect recovery depends primarily on the feature noise ρ , which can tolerate more with increasing nodes per clusters.

Case 3: Graph and node-specific information are synergistic. When $0 < \mu < n \cdot n_i$, both graph structure and node features contribute jointly. Two scenarios may arise:

1. If $\lambda \geq a \cdot p_{\max} n + \mu \zeta$, then combining with $\lambda \leq (1 - 2p_{\max}) n_i$ leads to

$$\rho \le \left(\frac{1 - (a \cdot n/n_i + 2) p_{\text{max}}}{\mu \ell} + \gamma\right) n_i, \tag{60}$$

$$p_{\max} \le \frac{1 - \mu \ell(\rho/n_i - \gamma)}{2 + a \cdot n/n_i}. \tag{61}$$

2. If $a \cdot p_{\max} n + \mu \zeta \ge \lambda \ge a \cdot p_{\max} n_i + \mu \zeta$, then combining with $\lambda \le (1 - 2p_{\max}) n_i$ gives

$$\left(\frac{1 - (a+2) p_{\max}}{\mu \ell} + \gamma\right) n_i \ge \rho \ge \left(\frac{1 - (a \cdot n/n_i + 2) p_{\max}}{\mu \ell} + \gamma\right) n_i, \quad (62)$$

$$\frac{1 - \mu \ell(\rho/n_i - \gamma)}{2 + a} \ge p_{\text{max}} \ge \frac{1 - \mu \ell(\rho/n_i - \gamma)}{2 + a \cdot n/n_i}.$$
 (63)

Combining both scenarios, the feasible region for perfect recovery can be summarized as

$$\rho \le \left(\frac{1 - (a+2) p_{\max}}{\mu \ell} + \gamma\right) n_i,\tag{64}$$

$$p_{\max} \le \frac{1 - \mu \ell(\rho/n_i - \gamma)}{2 + a}.\tag{65}$$

This highlights the advantage of combining both sources of information: perfect recovery is achievable under sparser graphs (smaller $p_{\rm max}$) and noisier features (larger ρ) than in the graph-only or node-only regimes.

B.5 FINAL THEOREM

Theorem B.1. Consider a fixed graph G with n nodes partitioned into K classes, where class i contains n_i nodes with $n_i \sim n_j$ for all i, j. Let ρ_{ij}^+ denote the relative misconnection rate between classes, and define

$$\rho^{+} := \max_{i,j} \rho_{ij}^{+}, \qquad p_{\max} = \rho^{+} + \delta, \qquad a = \max_{i,j} a_{ij}.$$

Fix $\gamma = \mu_1/\mu$, and suppose Assumptions 4.1–4.4 hold with fixed δ and R, and that the feasible set Θ is bounded. Then:

1. **Graph-only regime** ($\mu \to 0$). Perfect recovery is achievable if

$$p_{\max} \le \frac{1}{2 + a \cdot n/n_i}.$$

2. Node-only regime ($\mu = \Omega(n \cdot n_i)$). Perfect recovery is achievable if

$$\rho \leq \gamma n_i$$
.

3. Synergistic regime (0 < μ < $n \cdot n_i$). Perfect recovery is achievable if

$$\rho \le \left(\frac{1 - (a+2) p_{\text{max}}}{\mu \ell} + \gamma\right) n_i, \tag{66}$$

$$p_{\max} \le \frac{1 - \mu\ell(\rho/n_i - \gamma)}{2 + a}.\tag{67}$$

B.6 LEMMAS

Lemma B.2. Let $\theta_v := \sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,v}^2 \theta_a$ and $\mathbf{L} := \sum_{a \in \mathcal{A}} \lambda_a e \mathbf{e}^\top$, where each $\theta_a \in \Theta$ and $\|\mathbf{e}\| = 1$. Assume that for all $v \in \mathcal{V}$, $\sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,v}^2 = 1$. Then, it follows that $\theta_v \in \Theta$, $L_{uv} \leq 1$ for all $u, v \in \mathcal{V}$, and $\sum_{a \in \mathcal{A}} \lambda_a = n$.

Proof. If $\sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,v}^2 = 1$ holds, we have:

- Since each $\theta_a \in \Theta$ and $\sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,v}^2 = 1$, it follows that θ_v is a convex combination of elements in Θ , and hence $\theta_v \in \Theta$ due to the convexity of Θ .
- For $m{L} = \sum_{a \in \mathcal{A}} \lambda_a e m{e}^{ op}$, the Cauchy–Schwarz inequality implies

$$L_{uv} = \sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,u} \epsilon_{a,v} \le \sqrt{\sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,u}^2} \cdot \sqrt{\sum_{a \in \mathcal{A}} \lambda_a \epsilon_{a,v}^2} = 1.$$

• Summing the constraint over all v, considering $||e||^2 = 1$, gives:

$$\sum_{v} \sum_{a} \lambda_{a} \epsilon_{a,v}^{2} = \sum_{a} \lambda_{a} \sum_{v} \epsilon_{a,v}^{2} = n.$$

Lemma B.3. Consider the symmetric block matrix

$$\Gamma = \begin{bmatrix} \Gamma_{\mathbf{1}\mathbf{1}} & \Gamma_{\mathbf{1}\mathbf{2}} \\ \Gamma_{\mathbf{1}\mathbf{2}}^\top & \Gamma_{\mathbf{2}\mathbf{2}} \end{bmatrix},$$

where Γ_{11} , Γ_{22} are symmetric matrices satisfying $\Gamma_{11} \succeq \alpha I$ and $\Gamma_{22} \succeq \beta I$ for some $\alpha, \beta \geq 0$. Then D is positive semidefinite if $\sigma^2(\Gamma_{12}) \leq \alpha \beta$, where $\sigma(\Gamma_{12})$ denotes the largest singular value of B.

Proof. Take any vector $\mathbf{x} = [\mathbf{x}_1^\top, \mathbf{x}_2^\top]^\top$. Then, the quadratic form associated with Γ can be written as:

$$\mathbf{x}^{\mathsf{T}} \mathbf{\Gamma} \mathbf{x} = \mathbf{x}_{1}^{\mathsf{T}} \mathbf{\Gamma}_{11} \mathbf{x}_{1} + 2 \mathbf{x}_{1}^{\mathsf{T}} \mathbf{\Gamma}_{12} \mathbf{x}_{2} + \mathbf{x}_{2}^{\mathsf{T}} \mathbf{\Gamma}_{22} \mathbf{x}_{2}.$$

Using the given conditions $\mathbf{A} \succeq \alpha \mathbf{I}$ and $\mathbf{\Gamma_{22}} \succeq \beta \mathbf{I}$, we have: $\mathbf{x}_1^{\top} \mathbf{\Gamma_{11}} \mathbf{x}_1 \geq \alpha \|\mathbf{x}_1\|^2, \mathbf{x}_2^{\top} \mathbf{\Gamma_{22}} \mathbf{x}_2 \geq \beta \|\mathbf{x}_2\|^2$. Applying these lower bounds and the definition of singular value ($\|\mathbf{\Gamma_{11}} \mathbf{x}_2\| \leq \sigma(\mathbf{\Gamma_{12}}) \|\mathbf{x}_2\|$), it follows by the Cauchy–Schwarz inequality that:

$$\mathbf{x}^{\top} \mathbf{\Gamma} \mathbf{x} \ge \alpha \|\mathbf{x}_1\|^2 + \beta \|\mathbf{x}_2\|^2 - 2\sigma(\mathbf{\Gamma}_{12}) \|\mathbf{x}_1\| \|\mathbf{x}_2\|.$$

Now, complete the square explicitly to factorize the expression clearly:

$$\mathbf{x}^{\top} \mathbf{\Gamma} \mathbf{x} \ge (\sqrt{\alpha} \|\mathbf{x}_1\| - \sqrt{\beta} \|\mathbf{x}_2\|)^2 + 2(\sqrt{\alpha\beta} - \sigma(\mathbf{\Gamma}_{12})) \|\mathbf{x}_1\| \|\mathbf{x}_2\|.$$

Since $\sigma^2(\Gamma_{12}) \leq \alpha\beta$ (implying $\sigma(\Gamma_{12}) \leq \sqrt{\alpha\beta}$), both terms in this expression are nonnegative. Thus:

$$\mathbf{x}^{\top} \mathbf{\Gamma} \mathbf{x} > 0$$
, for all \mathbf{x} .

Hence, Γ is positive semi-definite.

Lemma B.4 (Spectral bound from sparsity and entry size). Let $\mathbf{M} \in \mathbb{R}^{n \times n}$ be symmetric. Suppose each row has at most d_{\max} nonzero entries and $|M_{ij}| \leq a$ for all i, j. Then

$$\lambda_{\max}(\mathbf{M}) = \|\mathbf{M}\|_2 \le d_{\max} a.$$

Proof. Since M is symmetric, $\lambda_{\max}(\mathbf{M}) = \|\mathbf{M}\|_2$. Also,

1136
1137
$$\|\mathbf{M}\|_{\infty} = \max_{i} \sum_{j=1}^{n} |M_{ij}| \le d_{\max} a.$$
1138

Because M is symmetric, $\|\mathbf{M}\|_1 = \|\mathbf{M}\|_{\infty}$. Hence

$$\|\mathbf{M}\|_{2} \le \sqrt{\|\mathbf{M}\|_{1} \|\mathbf{M}\|_{\infty}} = \|\mathbf{M}\|_{\infty} \le d_{\max} a.$$

1142 This yields the bound $\lambda_{\max}(\mathbf{M}) \leq d_{\max}a$.

Remark B.5. The bound $\lambda_{\max}(\mathbf{M}) \leq a \, d_{\max}$ is tight up to constants in general. For example, if \mathbf{M} is the adjacency matrix of a d-regular graph with a=1, then $\lambda_{\max}(\mathbf{M})=d=d_{\max}$.

Lemma B.6 (Lower bound for a projected principal block). Let $S \in \mathbb{R}^{n \times n}$ be symmetric. Fix any index set $B \subseteq \{1, ..., n\}$ and let $S_{B,B}$ denote the corresponding principal submatrix. Assume:

- 1. Each row of $S_{\mathcal{B},\mathcal{B}}$ has at most $d_{\max}(\mathcal{B})$ nonzeros;
- 2. $|S_{ij}| \leq a$ for all i, j.

Let \mathbf{P}^{\perp} be any orthogonal projector ($\mathbf{P}^{\perp^2} = \mathbf{P}^{\perp} = \mathbf{P}^{\perp^{\top}}$), and define

$$\mathbf{M}_{\mathcal{B}} := (\mathbf{P}^{\perp}(\lambda \mathbf{I} - \mathbf{S})\mathbf{P}^{\perp})_{\mathcal{B},\mathcal{B}}.$$

Then

$$\lambda_{\min}(\mathbf{M}_{\mathcal{B}}) \geq \min \Big\{ 0, \ \lambda - a \, d_{\max}(\mathcal{B}) \Big\}.$$

Proof. Since \mathbf{P}^{\perp} is an orthogonal projector, $\|\mathbf{P}^{\perp}\|_{2} = 1$ and

$$\mathbf{M}_{\mathcal{B}} = \left(\lambda \mathbf{P}^{\perp} - \mathbf{P}^{\perp} \mathbf{S} \mathbf{P}^{\perp}\right)_{\mathcal{B} \mathcal{B}}$$

Hence

$$\lambda_{\min}(\mathbf{M}_{\mathcal{B}}) \; \geq \; \min \Big\{ 0, \; \lambda_{\min} \left((\lambda \mathbf{I} - \mathbf{S})_{\mathcal{B}, \mathcal{B}} \right) \Big\} \; \geq \; \min \Big\{ 0, \; \lambda - \lambda_{\max} \left((\mathbf{S})_{\mathcal{B}, \mathcal{B}} \right) \Big\},$$

where we used $\|\mathbf{P}^{\perp}\|_2 = 1$ and submultiplicativity to drop \mathbf{P}^{\perp} in the norm bound on the $\mathcal{B} \times \mathcal{B}$ block.

By Lemma B.4 applied to $S_{\mathcal{B},\mathcal{B}}$,

$$\lambda_{\max}\left((\mathbf{S})_{\mathcal{B},\mathcal{B}} \right) \leq a \, d_{\max}(\mathcal{B}),$$

which yields

$$\lambda_{\min}(\mathbf{M}_{\mathcal{B}}) \geq \min \left\{ 0, \ \lambda - a \, d_{\max}(\mathcal{B}) \right\}.$$

C POLYNOMIAL-TIME CONVERGENCE

Our analysis builds on the concept of a *linear minimization oracle (LMO)*, a standard tool in conditional gradient methods. The LMO for a dictionary A is defined as follows:

Definition C.1 (LMO). The linear minimization oracle (LMO) of a dictionary \mathcal{A} is a map $\mathcal{O}_{\mathcal{A}}(\mathbf{w})$ that, for any vector \mathbf{w} , returns an atom $\mathbf{a} \in \mathcal{A}$ minimizing the inner product $\mathbf{a}^{\top}\mathbf{w}$.

To establish our convergence result, we make the following assumption:

Assumption C.2. The feasible set \mathcal{U} is closed and convex, the atomic set \mathcal{A} is bounded, and the function $\phi(U)$ is convex and lower semi-continuous. Moreover, there exists an optimal solution $(U^*, \{\lambda_i^*\}_i)$ of (5) and vectors z^* , \mathbf{n}^* , \mathbf{g}^* satisfying

$$\mathbf{n}^* \in \partial I_{\mathcal{U}}(\mathbf{U}^*), \quad \sup_{\mathbf{a} \in \mathcal{A}} \mathbf{a}^{\top} \mathbf{z}^* \le \mu_0, \quad \mathbf{g}^* \in \partial \phi(\mathbf{U}^*), \quad \mathbf{0} \in \mathbf{g}^* + \mathbf{z}^* + \mathbf{n}^*.$$
 (68)

Under this assumption, we obtain the following result:

Theorem C.3. Under Assumption C.2, there exists an algorithm that returns an ϵ -approximate solution of (5) in $\operatorname{poly}(1/\epsilon)$ number of LMOs, proximal evaluations of ϕ , and orthogonal projections onto \mathcal{U} .

The proof of Theorem C.3 is provided in the following. This result also extends to (6) by replacing ϕ with its composite form $\phi' := \phi + \mu_1 R$. While this establishes the polynomial-time solvability of the problem, the construction is primarily of theoretical interest and is not computationally efficient for large-scale instances. In particular, solving (5) exactly in practice is challenging due to the additional feasibility constraint $U \in \mathcal{U}$. Although algorithms such as CoGEnT (Rao et al., 2015) efficiently handle the unconstrained version of the problem, extending them to the constrained setting remains an open question. To address these practical challenges, we turn to a scalable approach that operates on the non-convex formulation of the problem with fixed number of atoms and provides reliable performance in practice.

C.1 Proof of Theorem C.3

For convenience, we first provide a brief overview of the proof techniques employed.

C.1.1 PROOF OVERVIEW

We cast the problem in (5) as a constrained atomic norm minimization, which we solve using a tailored ADMM-based algorithm. This method decouples the atomic norm regularization from additional structural constraints via a variable splitting strategy. To establish the complexity result, we demonstrate that our problem meets the conditions required for standard ADMM convergence, yielding an overall rate of O(1/T).

C.1.2 CONSTRAINED ATOMIC NORM OPTIMIZATION (CANO)

We consider the following general constrained atomic norm optimization problem, referred to as CANO:

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) \quad \text{subject to} \quad \|\boldsymbol{x}\|_{\mathcal{A}} \le \tau, \quad \boldsymbol{x} \in \mathcal{X}, \tag{69}$$

where $\|\cdot\|_{\mathcal{A}}$ is the atomic norm, and $\mathcal{X}\subseteq\mathbb{R}^d$ is a closed and convex set encoding additional constraints (e.g., feasibility constraints in (5)). This formulation introduces computational challenges due to the coupling of a non-polyhedral norm and convex constraint sets.

C.1.3 CANO-ADMM ALGORITHM

To solve (69), we apply ADMM by introducing an auxiliary variable z and rewriting the problem as:

$$\min_{x \in \mathcal{I}} f(x) \quad \text{subject to} \quad x = z, \quad \|x\|_{\mathcal{A}} \le \tau, \quad z \in \mathcal{X}. \tag{70}$$

The augmented Lagrangian is defined as:

$$\mathcal{L}_{\beta}(\boldsymbol{x}, \boldsymbol{z}, \boldsymbol{\lambda}) = f(\boldsymbol{x}) + \langle \boldsymbol{\lambda}, \boldsymbol{x} - \boldsymbol{z} \rangle + \frac{\beta}{2} \|\boldsymbol{x} - \boldsymbol{z}\|_{2}^{2}, \tag{71}$$

where λ is the dual variable, and $\beta>0$ is the penalty parameter. The ADMM updates proceed as follows:

• Update x:

$$\boldsymbol{x}_{t+1} = \arg\min_{\|\boldsymbol{x}\|_{\mathcal{A}} \leq \tau} f(\boldsymbol{x}) + \langle \boldsymbol{\lambda}_t, \boldsymbol{x} - \boldsymbol{z}_t \rangle + \frac{\beta}{2} \|\boldsymbol{x} - \boldsymbol{z}_t\|_2^2.$$
 (72)

This step is solved using the CoGEnT algorithm (Rao et al., 2015), which is designed for atomic norm-constrained problems.

• Update z:

$$\boldsymbol{z}_{t+1} = \arg\min_{\boldsymbol{z} \in \mathcal{X}} \langle \boldsymbol{\lambda}_t, \boldsymbol{x}_{t+1} - \boldsymbol{z} \rangle + \frac{\beta}{2} \| \boldsymbol{x}_{t+1} - \boldsymbol{z} \|_2^2.$$
 (73)

This corresponds to projecting onto the set \mathcal{X} , which is assumed to be tractable via a gradient-projection update:

$$z_{t+1} = P_{\mathcal{X}} (z_t + \alpha(\lambda_t + \beta(x_{t+1} - z_t))),$$

where α is a step size and $P_{\mathcal{X}}$ denotes Euclidean projection onto \mathcal{X} .

• Update λ :

$$\lambda_{t+1} = \lambda_t + \beta(\boldsymbol{x}_{t+1} - \boldsymbol{z}_{t+1}). \tag{74}$$

C.1.4 Convergence Guarantees of CANO-ADMM

The convergence of ADMM with a rate of O(1/T) in objective residuals and constraint violations is guaranteed under the following conditions:

- 1. f(x) is convex and has a Lipschitz-continuous gradient;
- 2. The constraint sets $||x||_{\mathcal{A}} \leq \tau$ and \mathcal{X} are both closed and convex;
- 3. The subproblems are solvable to sufficient accuracy at each iteration.

Our setup satisfies all these assumptions:

- The function f(x) is convex and differentiable (e.g., it is the sum of a linear term and convex node-wise losses);
- The atomic norm ball $\|x\|_{\mathcal{A}} \le \tau$ is convex by definition, and \mathcal{X} is assumed to be convex and compact;
- The projection and LMO steps in CoGEnT are computationally tractable and converge efficiently.

Hence, CANO-ADMM converges at a rate O(1/T) in the number of iterations. As each iteration requires a linear minimization oracle and projection operation, the overall runtime is polynomial in $1/\epsilon$, completing the proof of Theorem C.3.

D DETAILS OF THE CADO ALGORITHM

This section provides the full implementation details of the CADO algorithm. We describe how the embedding vectors and model parameters are updated using conditional gradient steps, relying on problem-specific linear minimization oracles (LMOs). Each component is addressed in a separate subsection, followed by the complete algorithm pseudo-code.

D.1 EMBEDDING UPDATE VIA CONDITIONAL GRADIENT

In this step, we update the embedding vectors $\{\bar{e}_i\}_{i=1}^r$, which define the low-rank matrix $L = \sum_i \bar{e}_i \bar{e}_i^{\mathsf{T}}$, while keeping the class-wise models $\{\bar{\theta}_i\}$ fixed. The update is performed by solving the following LMO:

$$\arg \min_{\{\bar{\boldsymbol{e}}_i\}_{i=1}^r} \quad \sum_{i=1}^r \langle \nabla_{\bar{\boldsymbol{e}}_i} \phi, \, \bar{\boldsymbol{e}}_i \rangle$$
s.t. $\boldsymbol{L} = \sum_i \bar{\boldsymbol{e}}_i \bar{\boldsymbol{e}}_i^\top \in \mathcal{B}, \quad \boldsymbol{\theta}_v = \sum_i \bar{\epsilon}_{i,v}^2 \bar{\boldsymbol{\theta}}_i \in \Theta \quad \forall v \in \mathcal{V},$ (75)

where ϕ is the objective function from (4):

$$\phi = -\langle \bar{\boldsymbol{A}}, \boldsymbol{L} \rangle + \mu \sum_{v \in \mathcal{V}} f_v \left(\boldsymbol{z}_v; \, \boldsymbol{\theta}_v \right).$$

Constraint simplification. We now simplify the constraints using the structure of L and θ_v . First, the constraint $L \in \mathcal{B}$ requires: $L_{vv} = 1$ for all $v, 0 \le L_{uv} \le 1$ for all u, v, and symmetry.

Given $L = \sum_i \bar{e}_i \bar{e}_i^{\top}$, the diagonal condition $L_{vv} = 1$ becomes:

$$L_{vv} = \sum_{i} \bar{\epsilon}_{i,v}^2 = 1.$$

This implies that the vector $\bar{\boldsymbol{\epsilon}}_v = (\bar{\epsilon}_{i,v})_i$ lies on the unit sphere, and that:

$$oldsymbol{ heta}_v = \sum_i ar{\epsilon}_{i,v}^2 ar{oldsymbol{ heta}}_i$$

is a convex combination of class-wise models, hence satisfying $\theta_v \in \Theta$ automatically. Also, under this condition:

$$L_{uv} = \sum_{i} \bar{\epsilon}_{i,u} \bar{\epsilon}_{i,v} \le \sqrt{\sum_{i} \bar{\epsilon}_{i,u}^{2}} \cdot \sqrt{\sum_{i} \bar{\epsilon}_{i,v}^{2}} = 1,$$

and the upper bound constraint $L_{uv} \leq 1$ is also satisfied. The only remaining constraint is the nonnegativity of L_{uv} , which is equivalent to:

$$L_{uv} = \sum_{i} \bar{\epsilon}_{i,u} \bar{\epsilon}_{i,v} \ge 0.$$

Reparameterization. To simplify the optimization, we define: $W_{v,i} := \bar{\epsilon}_{i,v}^2$. Let $\mathbf{W} \in \mathbb{R}^{n \times r}$ collect these squared embedding values. Under this reparameterization, the low-rank matrix becomes $\mathbf{L} = \sum_{i=1}^r \sqrt{W_{:,i}} \sqrt{W_{:,i}}^{\mathsf{T}}$, and the node-specific models satisfy $\boldsymbol{\theta}_v = \sum_{i=1}^r W_{v,i} \bar{\boldsymbol{\theta}}_i$. The constraint $L_{vv} = \sum_i W_{v,i} = 1$ together with $W_{v,i} \geq 0$ implies that each row of \mathbf{W} lies in the probability simplex.

Thus, the LMO reduces to:

$$\arg\min_{\mathbf{W}\in\mathbb{R}^{n\times r}} -\left\langle \bar{A}, \sum_{i=1}^{r} \sqrt{W_{:,i}} \sqrt{W_{:,i}}^{\top} \right\rangle + \mu \sum_{v\in\mathcal{V}} f_v \left(\mathbf{z}_v; \sum_{i=1}^{r} W_{v,i} \bar{\boldsymbol{\theta}}_i \right)$$
s.t.
$$\sum_{i} W_{v,i} = 1, \quad W_{v,i} \geq 0 \quad \forall v \in \mathcal{V}.$$

$$(76)$$

We omit the non-negativity constraint $L \ge 0$, as it is empirically satisfied due to the structure of \bar{A} . If needed, it can be enforced via ADMM.

Gradient Computation. Let $\mathbf{R}_v = \sum_i W_{v,i} \bar{\mathbf{R}}_i$, $\mathbf{\pi}_v = \sum_i W_{v,i} \bar{\mathbf{\pi}}_i$. The gradient of the objective in (76) with respect to $W_{v,i}$ consists of:

· Structural term:

$$\frac{\partial}{\partial W_{v,i}} \left(-\langle \bar{\boldsymbol{A}}, \boldsymbol{L} \rangle \right) = -\sum_{u \in \mathcal{V}} \bar{A}_{u,v} \frac{\sqrt{W_{u,i}}}{\sqrt{W_{v,i}}}$$

• Feature loss:

$$\frac{\partial}{\partial W_{v,i}} f_{\text{feature}}(\boldsymbol{x}_v; \boldsymbol{R}_v) = \frac{1}{m} \left(-\boldsymbol{x}_v^\top \boldsymbol{R}_v^{-1} \bar{\boldsymbol{R}}_i \boldsymbol{R}_v^{-1} \boldsymbol{x}_v + \text{Tr}(\bar{\boldsymbol{R}}_i) \right)$$

• Label loss (training nodes only):

$$\frac{\partial}{\partial W_{v,i}} f_{\text{label}}(\tilde{y}_v; \boldsymbol{\pi}_v) = -\bar{\pi}_{i, \tilde{y}_v}$$

Solution. The optimization in (76) is linear in each row $\mathbf{W}_{v,:}$ and constrained to the simplex. Therefore, the optimal solution is a one-hot vector with 1 at the coordinate corresponding to the minimum gradient value:

$$\tilde{W}_{v,i} = \begin{cases} 1, & i = \arg\min_{j} \nabla_{W_{v,j}} \phi, \\ 0, & \text{otherwise.} \end{cases}$$

$$\nabla_{W_{v,i}}\phi = \frac{\partial}{\partial W_{v,i}} \left(-\langle \bar{\boldsymbol{A}}, \boldsymbol{L} \rangle \right) + \mu \sum_{v \in \mathcal{V}} \frac{\partial}{\partial W_{v,i}} f_{\text{feature}}(\boldsymbol{x}_v; \boldsymbol{R}_v) + \mu \beta \sum_{v \in \mathcal{V}_{\text{train}}} \frac{\partial}{\partial W_{v,i}} f_{\text{label}}(\tilde{y}_v; \boldsymbol{\pi}_v)$$

Once W is computed, we recover embeddings via the following equations. Taking positive roots preserves the symmetry and ensures non-negativity of L.

$$\bar{\epsilon}_{i,v} = \sqrt{W_{v,i}}, \quad \bar{e}_i = (\bar{\epsilon}_{i,v})_{v=1}^n.$$

MODEL UPDATE VIA CONDITIONAL GRADIENT

Given the updated embedding vectors $\{\bar{e}_i\}_{i=1}^r$, we update the class-wise models $\{\bar{\theta}_i = (\bar{R}_i, \bar{\pi}_i)\}_{i=1}^r$ by solving one LMO over each atom. Specifically, each model parameter is updated via:

$$\arg\min_{\bar{\boldsymbol{\theta}}_i \in \Theta} \left\langle \nabla_{\bar{\boldsymbol{\theta}}_i} \phi, \, \bar{\boldsymbol{\theta}}_i \right\rangle, \qquad \forall i \in \{1, \dots, r\}, \tag{77}$$

where ϕ is the global objective from (4), and the embedding matrix **W** (with $W_{v,i} = \bar{\epsilon}_{i,v}^2$) is fixed.

Since each $\bar{\theta}_i$ consists of a feature model $\bar{R}_i \in \mathbb{S}_{\rho_-,\rho_+}$ and a label distribution $\bar{\pi}_i \in \Delta$, the LMO separates into two independent problems:

$$\arg \min_{\bar{\boldsymbol{R}}_{i} \in \mathbb{S}_{\rho_{-},\rho_{+}}} \left\langle \nabla_{\bar{\boldsymbol{R}}_{i}} \phi, \, \bar{\boldsymbol{R}}_{i} \right\rangle, \tag{78}$$

$$\arg \min_{\bar{\boldsymbol{\Pi}}_{i}} \left\langle \nabla_{\bar{\boldsymbol{\pi}}_{i}} \phi, \, \bar{\boldsymbol{\pi}}_{i} \right\rangle. \tag{79}$$

$$\arg \quad \min_{\bar{\boldsymbol{\pi}}_i \in \Delta} \quad \langle \nabla_{\bar{\boldsymbol{\pi}}_i} \phi, \, \bar{\boldsymbol{\pi}}_i \rangle \,. \tag{79}$$

Gradient computation. Let $R_v = \sum_{i=1}^r W_{v,i} \bar{R}_i$ and $\pi_v = \sum_{i=1}^r W_{v,i} \bar{\pi}_i$. Then, the gradients are given by:

• Feature loss (all nodes):

$$abla_{ar{m{R}}_i} \phi = \mu \sum_{v \in \mathcal{V}} W_{v,i} \cdot
abla_{m{R}_v} f_{ ext{feature}}(m{x}_v; m{R}_v),$$

where

$$abla_{m{R}_v} f_{ ext{feature}}(m{x}_v; m{R}_v) = rac{1}{m} \left(-m{R}_v^{-1} m{x}_v m{x}_v^{ op} m{R}_v^{-1} + \mathbf{I}
ight).$$

• Label loss (training nodes only):

$$\nabla_{\bar{\boldsymbol{\pi}}_i} \phi = \mu \beta \sum_{v \in \mathcal{V}_{\text{train}}} W_{v,i} \cdot \nabla_{\boldsymbol{\pi}_v} f_{\text{label}}(\tilde{y}_v; \boldsymbol{\pi}_v),$$

where

$$\nabla_{\boldsymbol{\pi}_v} f_{\text{label}}(\tilde{y}_v; \boldsymbol{\pi}_v) = -\boldsymbol{e}_{\tilde{y}_v}$$

Closed-form solutions. Both optimization problems admit closed-form solutions:

• Update for \bar{R}_i : The optimal solution to the linear minimization problem

$$\arg\min_{\bar{\boldsymbol{R}}_{i}\in\mathbb{S}_{\rho_{-},\rho_{+}}}\quad\left\langle \nabla_{\bar{\boldsymbol{R}}_{i}}\phi,\,\bar{\boldsymbol{R}}_{i}\right\rangle ,$$

is given by

$$\bar{\mathbf{R}}_i = U \operatorname{diag}(r_1, \dots, r_m) U^{\top},$$

where $\nabla_{\bar{R}_i} \phi = U \operatorname{diag}(\lambda_1, \dots, \lambda_m) U^{\top}$ is the eigen-decomposition of the gradient, and

$$r_k = \begin{cases} \rho_- & \text{if } \lambda_k > 0, \\ \rho_+ & \text{if } \lambda_k < 0, \\ \text{any value in } [\rho_-, \rho_+] & \text{if } \lambda_k = 0. \end{cases}$$

• Update for $\bar{\pi}_i$: Since the objective is linear over the simplex, the solution is the vertex corresponding to the smallest coordinate:

$$ilde{\pi}_i = e_{k^\star}, \quad ext{where} \quad k^\star = rg\min_k \left[
abla_{ar{\pi}_i} \phi
ight]_k.$$

These updates define the model step in each iteration of the CADO algorithm.

D.3 THE SPECIALIZED CADO ALGORITHM

We now summarize the complete CADO algorithm specialized for the node classification setting studied in this paper. This algorithm solves the constrained atomic decomposition problem in (4) using a conditional gradient approach, alternating between updating the embedding vectors $\{\bar{e}_i\}$ and the class-wise models $\{\bar{\theta}_i = (\bar{R}_i, \bar{\pi}_i)\}$.

Embedding step. In each iteration, the embedding update seeks a direction that reduces the global objective ϕ while maintaining feasibility. To make this step efficient, we reparameterize the squared embedding entries as $W_{v,i} = \bar{\epsilon}_{i,v}^2$, which allows us to enforce both the diagonal constraint on \boldsymbol{L} and the convexity condition on $\boldsymbol{\theta}_v$. The resulting LMO admits a closed-form solution: each row of \boldsymbol{W} is set to a one-hot vector in the direction of steepest descent.

Model step. Given the updated embeddings, the model parameters $\bar{\theta}_i = (\bar{R}_i, \bar{\pi}_i)$ are updated by solving LMOs over the model space $\Theta = \mathbb{S}_{\rho_-,\rho_+} \times \Delta$. The gradients of the global objective ϕ with respect to both components are derived in closed form based on the structure of the loss functions. These LMOs also admit simple solutions: the covariance matrix \bar{R}_i is updated by projecting the negative gradient onto the spectral box, while the label distribution $\bar{\pi}_i$ is updated by selecting the coordinate with the smallest gradient value.

Alternating optimization. The algorithm alternates between these two steps, using a step size $\gamma_t = \frac{2}{t+2}$ at iteration t to compute convex combinations of the previous and newly computed atoms. This ensures feasibility at all iterations and convergence under standard assumptions. The resulting procedure is efficient, scalable, and compatible with a wide range of feature and label models.

Final algorithm. The complete specialized version of the CADO algorithm is presented in Algorithm 2 below.

Algorithm 2 CADO Algorithm (Specialized for our studied case in section 6.1) 1: **Input:** Number of atoms r, graph \bar{A} , node data $\{z_v\}$, step size sequence $\{\gamma_t = 2/t + 2\}$ 2: Initialize $\{\bar{e}_i^{(0)} \in \mathcal{U}\}_{i=1}^r, \{\bar{\theta}_i^{(0)} = (\bar{R}_i^{(0)}, \bar{\pi}_i^{(0)}) \in \Theta\}_{i=1}^r$ 3: for $t = 0, 1, 2, \ldots$ until convergence do // Embedding Update via Conditional Gradient 4: Compute $W_{v,i}^{(t)} = \overline{\epsilon}_{i,v}^{(t)2}$, and evaluate $\boldsymbol{R}_{v}^{(t)}, \boldsymbol{\pi}_{v}^{(t)}$ Compute gradient $\nabla_{W_{v,i}} \phi$ 5: 6: Solve embedding LMO; set $\tilde{W}_{v,i}^{(t)} = 1$ at minimum coordinate, 0 elsewhere 7: Set $\tilde{\epsilon}_{i,v}^{(t)} = \sqrt{\tilde{W}_{v,i}^{(t)}}$, and form $\tilde{e}_i^{(t)} = (\tilde{\epsilon}_{i,v}^{(t)})_v$ Update: $\bar{e}_i^{(t+1)} = (1 - \gamma_t)\bar{e}_i^{(t)} + \gamma_t\tilde{e}_i^{(t)}$ 8: 9: 10: // Model Update via Conditional Gradient 11: Compute $\nabla_{\bar{R}_i} \phi$, $\nabla_{\bar{\pi}_i} \phi$ using formulas in Appendix D.2 Solve model LMO; set 12: $\tilde{\boldsymbol{R}}_{i}^{(t)} = \mathcal{P}_{\mathbb{S}_{\rho_{-},\rho_{+}}}\left(-\nabla_{\bar{\boldsymbol{R}}_{i}}\phi\right), \qquad \tilde{\boldsymbol{\pi}}_{i}^{(t)} = \boldsymbol{e}_{k^{\star}}, \quad k^{\star} = \arg\min_{k}\left[\nabla_{\bar{\boldsymbol{\pi}}_{i}}\phi\right]_{k}$ Update: 13: $\bar{R}_{i}^{(t+1)} = (1 - \gamma_{t})\bar{R}_{i}^{(t)} + \gamma_{t}\tilde{R}_{i}^{(t)}, \quad \bar{\pi}_{i}^{(t+1)} = (1 - \gamma_{t})\bar{\pi}_{i}^{(t)} + \gamma_{t}\tilde{\pi}_{i}^{(t)}$ 14: **Return:** $\{ar{e}_i^{(T)}\}, \{ar{R}_i^{(T)}, ar{\pi}_i^{(T)}\}$

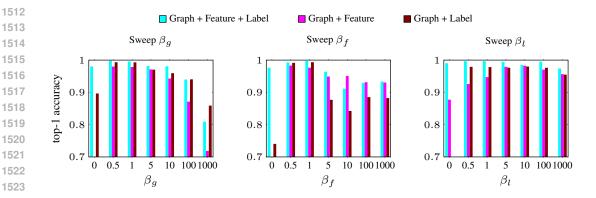


Figure 3: Sensitivity analysis of the proposed method with respect to the weighting parameters β_g , β_f , and β_l , which control the contribution of the graph structure, feature information, and label information, respectively, in the overall optimization objective. We report top-1 classification accuracy for three different configurations: Graph + Feature + Label (cyan), Graph + Feature (magenta), and Graph + Label (brown). Each plot varies one parameter while keeping the others fixed at their default values.

E EXTENDED EXPERIMENTS

E.1 MODEL PARAMETER SENSITIVITY

We study the sensitivity of the model's performance with respect to the weighting parameters β_g , β_f , and β_l , which control the influence of the graph term, feature term, and label term in our unified objective. The results are presented in Figure. 3.

In the left panel, we vary β_g , the weight of the graph term. When $\beta_g=0$, the model essentially ignores the graph structure, which causes a significant drop in performance in the Graph + Label configuration. However, the full model (Graph + Feature + Label) maintains high accuracy even at $\beta_g=0$, indicating that feature and label information alone can provide a strong signal. As β_g increases, accuracy improves across all settings, but plateaus after $\beta_g\approx 5$, suggesting diminishing returns from overly amplifying the graph signal.

In the center panel, we vary β_f , the weight of the feature term. When $\beta_f = 0$, the performance of the Graph + Feature configuration drops sharply, as expected. Interestingly, the full model remains relatively robust, highlighting the complementary strength of the graph and label terms. Very large values of β_f lead to performance degradation in some settings, possibly due to overfitting to noisy feature dimensions.

In the right panel, we sweep β_l , the weight of the label term. At $\beta_l = 0$, the Graph + Label configuration performs poorly due to the absence of label supervision. Again, the full model is quite robust and achieves high accuracy even with limited label contribution. Increasing β_l improves performance, but too high values lead to minor declines, likely because the model overemphasizes noisy or limited training labels.

Overall, these results confirm that our model is robust across a wide range of parameter choices and demonstrates strong synergy when all sources of information—graph, features, and labels—are integrated. The default values used in the main experiments strike a good balance across components.

E.2 DATA PARAMETER SENSITIVITY

Figure 4 investigates how the performance of our method changes as we vary key data generation parameters, specifically the number of nodes n, feature dimension m, and the number of clusters c.

In the left panel, increasing the number of nodes significantly improves test accuracy across all configurations, as larger graphs provide more structure and statistical signal. The full model (Graph + Feature + Label) consistently achieves the highest accuracy and converges quickly, even with a moderate number of nodes. This highlights the data efficiency of our joint framework, particularly when leveraging all available modalities.

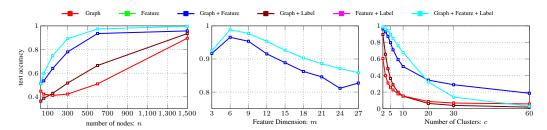


Figure 4: Effect of data parameters on node classification accuracy under different input configurations. Left: Accuracy versus total number of nodes n. Center: Accuracy versus total feature dimension m. Right: Accuracy versus number of clusters c. Each line corresponds to a different combination of information sources used: Graph, Feature, Label, or their combinations.

In the center panel, we vary the total feature dimension m, which includes both signal and noise components. As m increases, performance generally decreases—especially for methods that rely on features (e.g., Graph + Feature)—due to the increasing influence of noisy or uninformative dimensions. However, the full model (cyan) remains relatively stable and outperforms all other combinations, suggesting that combining features with graph structure and labels helps mitigate the curse of dimensionality.

In the right panel, we increase the number of clusters c, making the classification problem more challenging due to finer partitioning and weaker homophily. The performance of all configurations drops, but the full model retains significantly higher accuracy. Notably, using graph-only or label-only configurations fails beyond $c \approx 10$, while feature-based and joint models scale more gracefully. This result supports our theoretical findings: the integration of feature and label information compensates for reduced graph separability as the number of communities increases.

Overall, this experiment confirms the importance of combining multiple modalities. It also demonstrates the robustness of our approach across different data regimes, especially when individual signals become weak or insufficient.