

5.3 EXACT STRUCTURE RECOVERY VIA BOW PROJECTION

Finally, we formalize the post-hoc bow reconciliation (Section 4.3). For each pair $\{i, j\}$, define $d_{ij}^* := \max\{|B_{ij}^*|, |B_{ji}^*|\}$ and $r_{ij}^* := |\Omega_{ij}^*|/\sqrt{\Omega_{ii}^*\Omega_{jj}^*}$. Bow-freeness implies at most one of these is nonzero for each pair. We require a margin separating the two channels.

Assumption 5.3 (Pairwise margin). For fixed $c > 0$, define $\Delta_{ij}^* := |d_{ij}^* - cr_{ij}^*|$ and assume $\Delta^* := \min_{\{i,j\}} \Delta_{ij}^* > 0$.

Lemma 5.4 (Bow projection consistency). *Let $(\tilde{\mathbf{B}}, \tilde{\Omega})$ satisfy $\|\tilde{\mathbf{B}} - \mathbf{B}^*\|_\infty \leq \varepsilon_B$ and $\|\tilde{\Omega} - \Omega^*\|_\infty \leq \varepsilon_\Omega$. Under Assumption 5.3, if $\varepsilon_B + cL\varepsilon_\Omega < \frac{1}{2}\Delta^*$ for a constant $L > 0$ depending on the eigenvalue bounds, and post-hoc thresholds (τ_B, τ_Ω) are chosen appropriately, then the bow projection recovers the correct channel for every pair $\{i, j\}$, and the projected estimate $(\hat{\mathbf{B}}, \hat{\Omega})$ exactly recovers $\text{supp}(\mathbf{B}^*)$ and $\text{supp}_{\text{off}}(\Omega^*)$.*

Combining Theorem 5.2 with Lemma 5.4 yields the main structure recovery guarantee.

Corollary 5.5 (Exact support recovery). *Under the assumptions of Theorem 5.2 and Assumption 5.3, if the minimum signal strengths $\min_{(i,j) \in \text{supp}(\mathbf{B}^*)} |B_{ij}^*|$ and $\min_{(i,j) \in \text{supp}_{\text{off}}(\Omega^*)} |\Omega_{ij}^*|$ exceed $\sqrt{(\log p)/n}$ by a sufficient margin, then $\text{supp}(\hat{\mathbf{B}}) = \text{supp}(\mathbf{B}^*)$ and $\text{supp}_{\text{off}}(\hat{\Omega}) = \text{supp}_{\text{off}}(\Omega^*)$ with probability tending to one.*

6 EXPERIMENTS

We evaluate our method through comprehensive simulation studies and real-world datasets. The simulations systematically vary key structural parameters to assess identifiability and recovery performance across different regimes. As baselines, we compare against NOTEARS, GHOLE, GES, and DECAMF. DECAMF is designed to remove pervasive confounding effects and, in the linear setting, reduces to a two-step procedure: first removing a few principal components to eliminate low-rank latent structure, and then applying a structure learning method to the residualized data to estimate the sparse causal graph. For consistency and fair comparison, we employ NOTEARS in the second step.

We generate linear SEMs following the model in equation 1 with sparse directed edges \mathbf{B} and low-rank-plus-diagonal noise $\Omega = \mathbf{L}\mathbf{L}^\top + \sigma^2\mathbf{I}$. The generation process ensures bow-freeness through explicit cleanup: for any (i, j) pair where both $B_{ij} \neq 0$ and $\sum_k L_{ik}L_{jk} \neq 0$, we prioritize B_{ij} by zeroing out the common factor loadings in row j . For each configuration, we sample directed edges with $B_{ij} \sim \text{Uniform}([0.3, 0.8]) \times \text{sign}(\text{Rademacher})$ for randomly selected upper-triangular entries with density B_{density} , generate factor loadings where each column $\mathbf{L}_{:,k}$ has $\lfloor p \cdot L_{\text{density}} \rfloor$ non-zero entries drawn from $\mathcal{N}(0, 0.15^2)$, and generate data $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ where $\Sigma = (\mathbf{I} - \mathbf{B})^{-1}\Omega(\mathbf{I} - \mathbf{B})^{-\top}$. We set $\sigma^2 = 0.15$ throughout to maintain a consistent eigenvalue margin per Assumption ???. For each setting in each scenario, we generate 10 independent replicates. Unless specified otherwise, the sample complexity follows $n/p = 10$. We evaluate all methods on 10 independent replicates per density level, reporting mean performance with standard error bars.

We examine how latent confounding density affects causal structure recovery. We fix $p=20$ variables with structural density $B_{\text{density}}=0.1$, assume $q=5$ latent confounders, use $n=200$ samples, and vary $L_{\text{density}} \in \{0.0, 0.2, 0.4, 0.6, 0.8\}$ controlling the fraction of variable pairs influenced by shared latent factors. We compare DECOR_ADAPTIVE and DECOR_GL_ADAPTIVE (which adapt confounding regularization based on L_{density}) against ADAScore, NOTEARS, GOLEM, GES, LINGAM, and DECAMF.LIN variants. All methods use ℓ_1 penalty $\lambda_B=0.1$ for structure learning; adaptive methods adjust $\lambda_\Omega/\lambda_\Theta \in \{1.0, 0.5, 0.25, 0.1, 0.05\}$ corresponding to increasing L_{density} .

Performance Analysis. Figure 2 reveals several critical insights. First, *adaptive joint modeling dramatically outperforms sequential deconfounding*: DECOR_ADAPTIVE and ADAScore achieve remarkably low SHD (20–25) and FPR (~ 0.01) across all densities, while DECAMF.LIN_{rTrue} exhibits catastrophic failure with near-zero TPR and F1 scores below 0.1. This 3–6 \times SHD advantage confirms that joint estimation of \mathbf{B} and Ω is fundamentally superior to two-stage factor-analytic residualization, which destroys causal signal.

Second, *DECOR_GL_ADAPTIVE exhibits density-dependent instability*: while achieving competitive performance at moderate densities ($L_{\text{density}}=0.2\text{--}0.4$) with $F1 \approx 0.3$, it suffers catastrophic

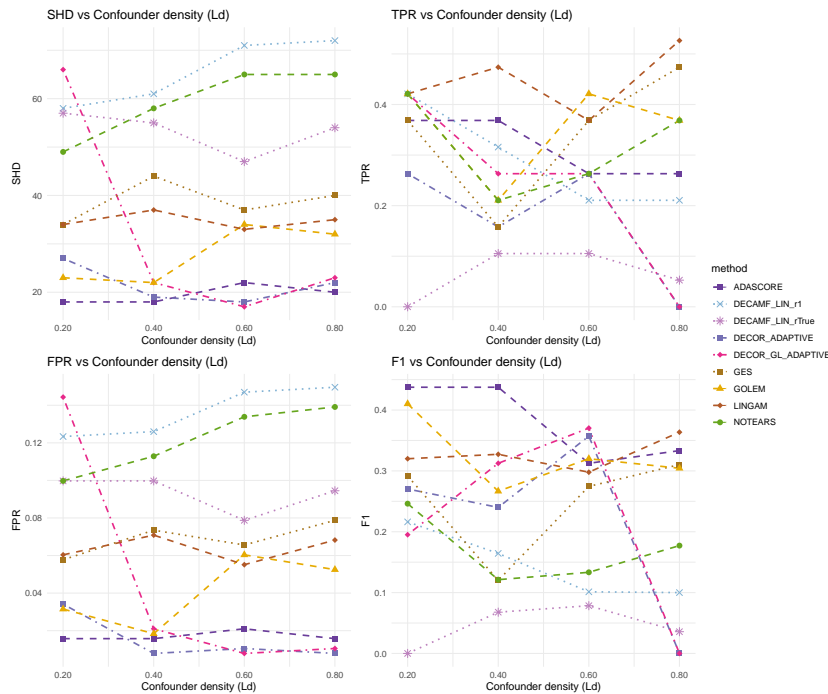


Figure 2: Performance under varying confounding density ($p=20$, $q=5$, $n=200$, $B_{\text{density}}=0.1$). DECOR_ADAPTIVE and ADASCOPE achieve 3–6 \times lower SHD than traditional methods; DECOR_GL_ADAPTIVE shows instability at extreme densities. Each curve shows mean across 10 replicates; error bars indicate standard errors.

degradation at high density ($L_{\text{density}}=0.8$), where TPR drops to nearly zero and F1 collapses. The U-shaped SHD curve (starting at 125 for $L_{\text{density}}=0.0$, dropping to 20–30 at moderate densities, maintaining 30 at high density) suggests the graphical lasso approach struggles in both unconfounded and heavily confounded regimes, likely due to ill-conditioning of the precision matrix Θ when confounding structure becomes dense.

Third, *traditional methods show graceful but significant degradation*: NOTEARS, GOLEM, GES, and LINGAM maintain relatively stable SHD (50–70) and moderate TPR (0.3–0.4) as confounding increases, but their consistently higher FPR (0.1–0.15) and lower F1 (0.1–0.3) demonstrate the cost of ignoring latent confounding. These methods still recover a meaningful subset of true edges but incur substantial false discoveries, confirming theoretical predictions that unmodeled confounders induce spurious conditional dependencies.

Fourth, the *precision-recall tradeoff* distinguishes method classes: DECOR_ADAPTIVE achieves optimal balance with moderate TPR (~ 0.25 – 0.27) but exceptionally low FPR (~ 0.01), yielding highest F1 (~ 0.45). Traditional methods exhibit inverted tradeoffs—higher TPR but 10 \times higher FPR—reflecting liberal edge declaration when confounding creates spurious correlations. The variance across replicates (error bars) is notably lower for DECOR_ADAPTIVE than constraint-based methods (GES), reflecting continuous optimization stability.

REFERENCES

- Raj Agrawal, Chandler Squires, Neha Prasad, and Caroline Uhler. The decamfounder: Non-linear causal discovery in the presence of hidden variables. *arXiv preprint arXiv:2102.07921*, 2023.
- Steen A Andersson, David Madigan, and Michael D Perlman. A characterization of markov equivalence classes for acyclic digraphs. *The Annals of Statistics*, 25(2):505–541, 1997.
- Hedy Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the kurdyka-

- 540 lojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.
- 541
- 542 Hedy Attouch, Jérôme Bolte, and Benar Fux Svaiter. Convergence of descent methods for semi-
543 algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized
544 Gauss–Seidel methods. *Mathematical Programming*, 137(1-2):91–129, 2013.
- 545 Kevin Bello, Bryon Aragam, and Pradeep K Ravikumar. Dagma: Learning dags via m-matrices
546 and a log-determinant acyclicity characterization. In *Advances in Neural Information Processing*
547 *Systems (NeurIPS)*, 2022.
- 548 Daniel Bernstein, Basil Saeed, Johannes Brehmer, Antti Hyttinen, and Caroline Uhler. Ordering-
549 based causal structure learning in the presence of latent variables (gspo). In *Advances in Neural*
550 *Information Processing Systems (NeurIPS)*, 2020.
- 551 Jérôme Bolte, Aris Daniilidis, and Adrian Lewis. The lojasiewicz inequality for nonsmooth suban-
552alytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimiza-*
553 *tion*, 17(4):1205–1223, 2007.
- 554
- 555 Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization
556 for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1–2):459–494, 2014.
- 557
- 558 Philippe Brouillard, Sébastien Lachapelle, Alexandre Lacoste, Simon Lacoste-Julien, and Alexandre
559 Drouin. Differentiable causal discovery from interventional data. *Advances in Neural Information*
560 *Processing Systems*, 33:21865–21877, 2020.
- 561 Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and*
562 *applications*. Springer Science & Business Media, 2011.
- 563
- 564 Venkat Chandrasekaran, Pablo A Parrilo, and Alan S Willsky. Latent variable graphical model
565 selection via convex optimization. In *2010 48th Annual Allerton Conference on Communication,*
566 *Control, and Computing (Allerton)*, pp. 1610–1613. IEEE, 2010.
- 567 David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of Machine*
568 *Learning Research*, 3:507–554, 2002a.
- 569 David Maxwell Chickering. Learning equivalence classes of bayesian-network structures. *Journal*
570 *of machine learning research*, 2(Feb):445–498, 2002b.
- 571
- 572 Diego Colombo, Marloes H. Maathuis, Markus Kalisch, and Thomas S. Richardson. Learning high-
573 dimensional directed acyclic graphs with latent and selection variables. *The Annals of Statistics*,
574 40(1):294–321, 2012.
- 575 Chang Deng, Kevin Bello, Pradeep Ravikumar, and Bryon Aragam. Markov equivalence and consis-
576 tency in differentiable structure learning. In *Advances in Neural Information Processing Systems*
577 *(NeurIPS)*, 2024.
- 578 Mathias Drton, Rina Foygel, and Seth Sullivant. Global identifiability of linear structural equation
579 models. 2011.
- 580
- 581 Bertrand Frot, Sven Nelander, and Caroline Uhler. Robust causal structure learning in the presence
582 of pervasive confounding. *arXiv preprint arXiv:1902.09057*, 2019.
- 583
- 584 Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear
585 causal discovery with additive noise models. *Advances in neural information processing systems*,
586 21, 2008.
- 587 Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based
588 neural dag learning. *arXiv preprint arXiv:1906.02226*, 2019.
- 589
- 590 Po-Ling Loh and Peter Bühlmann. High-dimensional learning of linear causal networks via inverse
591 covariance estimation. *Journal of Machine Learning Research*, 15(140):3065–3105, 2014.
- 592 Po-Ling Loh and Martin J. Wainwright. Regularized m-estimators with nonconvexity: Statistical
593 and algorithmic theory for local optima. *Journal of Machine Learning Research*, 16:559–616,
2015.

- 594 Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and dag constraints
595 for learning linear dags. *Advances in Neural Information Processing Systems*, 33:17943–17954,
596 2020.
- 597 Christopher Nowzohour, Marloes H. Maathuis, Robin J. Evans, and Peter Bühlmann. Distributional
598 equivalence and structure learning for bow-free acyclic path diagrams. *Electronic Journal of
599 Statistics*, 11(2):5342–5374, 2017. doi: 10.1214/17-EJS1372.
- 600 Samhita Pal, Dhruvajyoti Ghosh, and Shu Yang. Penalized fci for causal structure learning in a
601 sparse dag for biomarker discovery in parkinson’s disease. *arXiv preprint arXiv:2507.00173*,
602 2025.
- 603 Judea Pearl. *Causality*. Cambridge university press, 2009.
- 604 Jonas Peters and Peter Bühlmann. Identifiability of gaussian structural equation models with equal
605 error variances. *Biometrika*, 101(1):219–228, 2014.
- 606 Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Minimax rates of estimation for high-
607 dimensional linear regression over ℓ_q -balls. *IEEE Transactions on Information Theory*, 57(10):
608 6976–6994, 2011.
- 609 Agnieszka Reisach, Christoph Seiler, and Sebastian Weichwald. Beware of the simulated dag!
610 causal discovery benchmarks may be easy to game. *Advances in Neural Information Processing
611 Systems (NeurIPS)*, 2021.
- 612 Thomas Richardson and Peter Spirtes. Ancestral graph markov models. *The Annals of Statistics*, 30
613 (4):962–1030, 2002.
- 614 Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. Causal
615 protein-signaling networks derived from multiparameter single-cell data. *Science*, 308(5721):
616 523–529, 2005.
- 617 Andrea Seng, Ananya Ghosh, Steve Hanneke, and Bryon Aragam. Harder than you think: Consis-
618 tency of continuous optimization approaches for causal discovery. In *International Conference
619 on Learning Representations (ICLR)*, 2023.
- 620 Parikshit Shah, Jonas Peters, and Peter Bühlmann. Spectral deconfounding for causal structure
621 learning in linear models. In *Advances in Neural Information Processing Systems (NeurIPS)*,
622 2020.
- 623 Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. A linear
624 non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(10),
625 2006.
- 626 Kirankumar Shiragur, Jiaqi Zhang, and Caroline Uhler. Causal discovery with fewer conditional
627 independence tests. *arXiv preprint arXiv:2406.01823*, 2024.
- 628 Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT Press,
629 2nd edition, 2000a.
- 630 Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT
631 press, 2000b.
- 632 Chandler Squires and Caroline Uhler. Causal structure learning: A combinatorial perspective. *Founda-
633 tions of Computational Mathematics*, 23(5):1781–1815, 2023.
- 634 Chandler Squires, Annie Yun, Eshaan Nichani, Raj Agrawal, and Caroline Uhler. Causal structure
635 discovery between clusters of nodes induced by latent factors. In *Proceedings of the 25th Interna-
636 tional Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 177 of *Proceedings
637 of Machine Learning Research*, pp. 5267–5291, 2022. URL [https://proceedings.mlr.
638 press/v177/squires22a/squires22a.pdf](https://proceedings.mlr.press/v177/squires22a/squires22a.pdf).
- 639 Thijs van Ommen and Joris M. Mooij. Algebraic equivalence of linear structural equation models. In
640 *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence (UAI)*. AUAI Press,
641 2017.

- 648 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*.
649 Cambridge University Press, 2018.
- 650 Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge Uni-
651 versity Press, 2019.
- 652 Yibei Wang and Mathias Drton. Causal discovery with bow-free acyclic non-gaussian graphs.
653 *Journal of Machine Learning Research*, 24(315):1–45, 2023. URL [https://jmlr.org/
654 papers/v24/23-0217.html](https://jmlr.org/papers/v24/23-0217.html).
- 655 Yixin Wang and David M Blei. The blessings of multiple causes. *Journal of the American Statistical
656 Association*, 114(528):1574–1596, 2019.
- 657 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural
658 networks. In *International Conference on Machine Learning (ICML)*, 2019.
- 659 Jiji Zhang. Causal reasoning with ancestral graphs. *Journal of Machine Learning Research*, 9(7),
660 2008.
- 661 Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Contin-
662 uous optimization for structure learning. In *Advances in Neural Information Processing Systems
663 (NeurIPS)*, 2018.
- 664 Xun Zheng, Chen Dan, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Learning sparse
665 nonparametric dags. In *International Conference on Artificial Intelligence and Statistics (AIS-
666 TATS)*, pp. 3414–3425. Pmlr, 2020.

671 A PARAMETER IDENTIFIABILITY THEORY

672
673 *Proof of Lemma 2.3.* Let the variables be topologically ordered so that \mathbf{B} is strictly upper triangular
674 and $\mathbf{T} = (\mathbf{I} - \mathbf{B})^{-1}$ is unit lower triangular. For a node i , write $[i] = \{1, \dots, i\}$, parent set
675 $\text{Pa}(i) \subseteq [i - 1]$, sibling set $\text{Sib}(i) \subseteq [i - 1]$, and let

$$676 A_i := \Omega_{[i] \setminus \text{Sib}(i), [i]}, \quad B_i := \mathbf{T}_{[i], \text{Pa}(i)}.$$

677 The rank test at node i is that $A_i B_i$ has column rank $|\text{Pa}(i)|$.

678 Since \mathbf{B} is strictly upper triangular in a topological order, $\mathbf{T} = (\mathbf{I} - \mathbf{B})^{-1}$ is unit lower triangular.
679 Hence, for any i and any parent $j \in \text{Pa}(i) \subseteq [i - 1]$, the j -th row of $\mathbf{T}_{[i], \text{Pa}(i)}$ has a 1 in column
680 j and zeros in columns $\text{Pa}(i) \cap \{1, \dots, j - 1\}$. In particular, the row-selector R_i that keeps rows
681 $\text{Pa}(i)$ satisfies

$$682 R_i \mathbf{T}_{[i], \text{Pa}(i)} = I_{|\text{Pa}(i)|}.$$

683 Thus, for all $x \in \mathbb{R}^{|\text{Pa}(i)|}$, $\|\mathbf{T}_{[i], \text{Pa}(i)} x\| \geq \|R_i \mathbf{T}_{[i], \text{Pa}(i)} x\| = \|x\|$. Therefore $\sigma_{\min}(\mathbf{T}_{[i], \text{Pa}(i)}) \geq 1$,
684 and $\mathbf{T}_{[i], \text{Pa}(i)}$ has full column rank. \square

685
686
687 *Proof of Lemma 2.4.* Let $J := [i] \setminus \text{Sib}(i)$ and $A_i := \Omega_{J, [i]}$. We use two standard facts from
688 linear algebra: (i) every principal submatrix of a positive definite matrix is positive definite, and (ii)
689 eigenvalue interlacing implies that if $M \succeq m\mathbf{I}$, then every principal submatrix $M_{J, J}$ also satisfies
690 $M_{J, J} \succeq m\mathbf{I}$.

691 Since $\Omega \succeq m\mathbf{I}$ by Assumption 2.7, the principal block $\Omega_{[i], [i]}$ is positive definite with
692 $\lambda_{\min}(\Omega_{[i], [i]}) \geq m$. Consequently, $\Omega_{[i], [i]}^2$ is also positive definite with $\lambda_{\min}(\Omega_{[i], [i]}^2) =$
693 $\lambda_{\min}(\Omega_{[i], [i]})^2 \geq m^2$.

694 The Gram matrix of A_i can be written as

$$695 A_i A_i^\top = \Omega_{J, [i]} \Omega_{[i], J} = (\Omega_{[i], [i]}^2)_{J, J},$$

696 where the second equality uses symmetry of Ω . Since $A_i A_i^\top$ is a principal submatrix of $\Omega_{[i], [i]}^2$,
697 eigenvalue interlacing gives $\lambda_{\min}(A_i A_i^\top) \geq m^2$. Therefore,

$$698 \sigma_{\min}(A_i) = \sqrt{\lambda_{\min}(A_i A_i^\top)} \geq m,$$

699 and in particular A_i has full row rank. \square

Proof of Theorem 2.5. Fix i . By Lemma 2.3, $\sigma_{\min}(B_i) \geq 1$ and B_i has $|\text{Pa}(i)|$ independent columns. By Lemma 2.4, $\sigma_{\min}(A_i) \geq m > 0$, so A_i has full row rank. The bow-freeness assumption implies, $\text{Pa}(i) \cap \text{Sib}(i) = \emptyset$, hence the number of rows of A_i satisfies $|\text{Pa}(i)| \geq |\text{Pa}(i)|$, so the product $A_i B_i$ can (and will) have full column rank. Using the singular-value inequality again,

$$\sigma_{\min}(A_i B_i) \geq \sigma_{\min}(A_i) \sigma_{\min}(B_i) \geq m,$$

which implies $\text{rank}(A_i B_i) = |\text{Pa}(i)|$. Thus the node-wise rank condition holds for this i ; since i was arbitrary, it holds for all nodes. By the equivalence for acyclic graphs, the parametrization $(\mathbf{B}, \mathbf{\Omega}) \mapsto \Sigma$ is injective. \square

B OPTIMIZATION THEORY: ASSUMPTIONS AND PROOFS

This appendix states the assumptions underlying Section 5 and sketches the main proofs. We write the population negative log-likelihood as

$$\ell(\mathbf{B}, \mathbf{\Omega}) := -\Sigma^* [\log p_{\mathbf{B}, \mathbf{\Omega}}(X)],$$

with $\Sigma^* = \Sigma(\mathbf{B}^*, \mathbf{\Omega}^*)$, and the sample version as \mathcal{L}_n in equation 2. The penalized sample objective is \mathcal{J}_n in equation ??.

B.1 OBJECTIVE, PENALTIES, AND THE KL PROPERTY

We first spell out the penalty class and the KL structure.

Assumption B.1 (Amenable penalties). The penalties P_B and P_Ω are separable:

$$P_B(\mathbf{B}) = \sum_{i,j} p_B(|B_{ij}|), \quad P_\Omega(\mathbf{\Omega}) = \sum_{i \neq j} p_\Omega(|\Omega_{ij}|),$$

where $p_B, p_\Omega : [0, \infty) \rightarrow [0, \infty)$ are *amenable regularizers* in the sense of Loh & Wainwright (2015): they are continuous, differentiable on $(0, \infty)$, satisfy $p(0) = 0$, are nondecreasing and concave, and obey $|p'(t)| \leq \lambda$ and $p''(t) \geq -\mu$ for some (λ, μ) with μ small relative to the RSC constants. Both the ℓ_1 penalty ($p(t) = \lambda t$) and quasi-MCP satisfy these conditions. We further assume p_B and p_Ω are semi-algebraic.

Under Assumption B.1, P_B and P_Ω are proper, lower semicontinuous, and proximal. Moreover, they are semi-algebraic, hence KL functions. The smooth part

$$f_n(\mathbf{B}, \mathbf{\Omega}) := \mathcal{L}_n(\mathbf{B}, \mathbf{\Omega}) + \rho h(\mathbf{B})$$

is real-analytic on the domain $\{\mathbf{B} : \rho(|\mathbf{B}|) < 1\} \times \{\mathbf{\Omega} \succ 0\}$: the covariance map $(\mathbf{B}, \mathbf{\Omega}) \mapsto \Sigma(\mathbf{B}, \mathbf{\Omega})$ is analytic whenever $\mathbf{I} - \mathbf{B}$ is invertible, the Gaussian log-likelihood is analytic on the positive definite cone, and $h(\mathbf{B}) = \text{tr}(\exp(\mathbf{B} \circ \mathbf{B})) - p$ is analytic as a composition of polynomials, matrix exponential, and trace. Therefore $\mathcal{J}_n = f_n + P_B + P_\Omega$ is a KL function (Attouch et al., 2010; 2013; Bolte et al., 2007).

We further assume:

Assumption B.2 (Blockwise Lipschitz gradients and bounded level sets). On the set $\{\mathbf{\Omega} \succeq \eta \mathbf{I}\}$, the gradients $\nabla_{\mathbf{B}} f_n$ and $\nabla_{\mathbf{\Omega}} f_n$ are Lipschitz in each block, and the sublevel sets $\{(\mathbf{B}, \mathbf{\Omega}) : \mathcal{J}_n(\mathbf{B}, \mathbf{\Omega}) \leq c\}$ are bounded.

The Lipschitz property follows from standard matrix calculus and the eigenvalue margin on $\mathbf{\Omega}$; boundedness of level sets uses the facts that $\mathcal{L}_n \rightarrow \infty$ as $\mathbf{\Omega}$ approaches singularity or as \mathbf{B} approaches an unstable matrix, and that the penalties penalize large entries.

Proof of Proposition 5.1. Under Assumptions B.1–B.2, the DECOR updates coincide with one step of the PALM algorithm of Bolte et al. (2014): in Stage 1 we apply a proximal-gradient step in the \mathbf{B} -block with backtracking to ensure a descent inequality, and in Stage 2 we compute a proximal minimizer in the $\mathbf{\Omega}$ -block (or a proximal-gradient step, depending on the implementation), followed by an SPD projection that preserves boundedness of level sets. The PALM convergence theorem (Bolte et al., 2014, Theorem 3.1), specialized to KL objectives, implies: (i) monotone decrease of \mathcal{J}_n ; (ii) finite length of the iterates $\sum_t \|(\mathbf{B}^{(t+1)}, \mathbf{\Omega}^{(t+1)}) - (\mathbf{B}^{(t)}, \mathbf{\Omega}^{(t)})\| < \infty$; and (iii) convergence of the entire sequence to a single critical point of \mathcal{J}_n . This yields all claims in Proposition 5.1. \square

B.2 LOCAL CURVATURE AND STATISTICAL ERROR

We now formalize the local assumptions used in Theorem 5.2 and sketch the contraction argument.

Assumption B.3 (Local RSC and cross-Lipschitz). There exists a neighborhood \mathcal{N} of (\mathbf{B}^*, Ω^*) and constants $\mu_B, \mu_\Omega > 0$ and $L_{B\Omega}, L_{\Omega B} > 0$ such that for all $(\mathbf{B}, \Omega) \in \mathcal{N}$,

$$\begin{aligned} \langle \nabla_{\mathbf{B}} \ell(\mathbf{B}, \Omega) - \nabla_{\mathbf{B}} \ell(\mathbf{B}^*, \Omega), \mathbf{B} - \mathbf{B}^* \rangle &\geq \mu_B \|\mathbf{B} - \mathbf{B}^*\|_F^2, \\ \langle \nabla_{\Omega} \ell(\mathbf{B}, \Omega) - \nabla_{\Omega} \ell(\mathbf{B}, \Omega^*), \Omega - \Omega^* \rangle &\geq \mu_\Omega \|\Omega - \Omega^*\|_F^2, \end{aligned}$$

and

$$\begin{aligned} \|\nabla_{\mathbf{B}} \ell(\mathbf{B}, \Omega) - \nabla_{\mathbf{B}} \ell(\mathbf{B}, \Omega^*)\|_F &\leq L_{B\Omega} \|\Omega - \Omega^*\|_F, \\ \|\nabla_{\Omega} \ell(\mathbf{B}, \Omega) - \nabla_{\Omega} \ell(\mathbf{B}^*, \Omega)\|_F &\leq L_{\Omega B} \|\mathbf{B} - \mathbf{B}^*\|_F. \end{aligned}$$

Assumption B.4 (Sample-level deviations). The sample loss \mathcal{L}_n satisfies the same RSC and cross-Lipschitz bounds as ℓ , with constants shrunk by at most a factor of $1/2$, and $\|\nabla \mathcal{L}_n(\mathbf{B}^*, \Omega^*)\|_\infty = O(\sqrt{(\log p)/n})$.

Under sub-Gaussian tails, Assumption B.4 follows from standard matrix concentration inequalities (e.g. Vershynin, 2018; Wainwright, 2019). Bow-freeness and the eigenvalue margin imply Assumption B.3 by ensuring that the relevant population Hessians are well-conditioned; see Section 2.

Finally, we adopt the usual sparsity and signal-strength assumption:

Assumption B.5 (Sparsity and signal strength). Let $S_B = \text{supp}(\mathbf{B}^*)$ and $S_\Omega = \text{supp}_{\text{off}}(\Omega^*)$ with sizes s_B, s_Ω . For penalties with tuning parameters $\lambda_B, \lambda_\Omega \asymp \sqrt{(\log p)/n}$, assume

$$\min_{(i,j) \in S_B} |B_{ij}^*| \geq c_B \lambda_B, \quad \min_{(i,j) \in S_\Omega} |\Omega_{ij}^*| \geq c_\Omega \lambda_\Omega$$

for constants $c_B, c_\Omega > 1$.

Assumption B.5 is standard in high-dimensional sparse estimation; see, e.g., Bühlmann & Van De Geer (2011); Raskutti et al. (2011).

Proof of Theorem 5.2. Let $e_B^{(t)} = \|\mathbf{B}^{(t)} - \mathbf{B}^*\|_F$ and $e_\Omega^{(t)} = \|\Omega^{(t)} - \Omega^*\|_F$. The \mathbf{B} -update in DECOR minimizes (approximately) the function $\mathbf{B} \mapsto \mathcal{L}_n(\mathbf{B}, \Omega^{(t)}) + P_B(\mathbf{B}) + \rho h(\mathbf{B})$, keeping $\Omega^{(t)}$ fixed. By Assumptions B.3–B.4 and the amenable-regularizer analysis of Loh & Wainwright (2015), we obtain the one-step bound

$$e_B^{(t+1)} \leq \alpha_B e_B^{(t)} + \beta_B e_\Omega^{(t)} + C_B \sqrt{\frac{s_B \log p}{n}},$$

for some $\alpha_B < 1, \beta_B > 0$, and $C_B > 0$. The first term is a contraction from local strong convexity in \mathbf{B} , the second term captures the dependence on $\Omega^{(t)}$ via the cross-Lipschitz constants, and the last term is statistical error.

Similarly, the Ω -update minimizes $\Omega \mapsto \mathcal{L}_n(\mathbf{B}^{(t+1)}, \Omega) + P_\Omega(\Omega)$ for fixed $\mathbf{B}^{(t+1)}$. Since this subproblem is convex and locally strongly convex near Ω^* , standard arguments for ℓ_1 - or amenable-regularized covariance estimation yield

$$e_\Omega^{(t+1)} \leq \alpha_\Omega e_\Omega^{(t)} + \beta_\Omega e_B^{(t+1)} + C_\Omega \sqrt{\frac{s_\Omega \log p}{n}},$$

with $\alpha_\Omega < 1, \beta_\Omega > 0$, and $C_\Omega > 0$.

Combining the two inequalities and stacking $e^{(t)} = (e_B^{(t)}, e_\Omega^{(t)})^\top$, we obtain a linear recursion

$$e^{(t+1)} \leq M e^{(t)} + b_n,$$

for a 2×2 matrix M with spectral radius $\rho(M) < 1$ when \mathcal{N} is chosen small enough so that RSC dominates cross-Lipschitz effects, and b_n of order $\delta_n \asymp \sqrt{(s_B \log p)/n} + \sqrt{(s_\Omega \log p)/n}$. Solving the recursion yields $\|e^{(t)}\|_1 \leq \rho(M)^t \|e^{(0)}\|_1 + \|b_n\|_1 / (1 - \rho(M))$. This proves the claimed contraction and statistical error bound. \square

B.3 BOW PROJECTION CONSISTENCY

We expand Lemma 5.4. Recall from Section 5 that for each pair $\{i, j\}$,

$$\tilde{d}_{ij} = \max\{|\tilde{B}_{ij}|, |\tilde{B}_{ji}|\}, \quad \tilde{r}_{ij} = \frac{|\tilde{\Omega}_{ij}|}{\sqrt{\tilde{\Omega}_{ii}\tilde{\Omega}_{jj}}},$$

and similarly d_{ij}^*, r_{ij}^* for the true parameters.

Lemma B.6 (Lipschitz control of r_{ij}). *Suppose Ω^* satisfies $\eta\mathbf{I} \preceq \Omega^* \preceq M\mathbf{I}$ for some $0 < \eta \leq M$. There exists $L = L(\eta, M)$ such that for any positive definite $\tilde{\Omega}$ with $\|\tilde{\Omega} - \Omega^*\|_\infty \leq 1$,*

$$|\tilde{r}_{ij} - r_{ij}^*| \leq L \|\tilde{\Omega} - \Omega^*\|_\infty \quad \text{for all } \{i, j\}.$$

Proof. On the compact set $\mathcal{C} = \{\Omega : \frac{\eta}{2} \leq \Omega_{kk} \leq 2M, \|\Omega\|_\infty \leq 2M\}$ the map $r_{ij}(\Omega) = |\Omega_{ij}|/\sqrt{\Omega_{ii}\Omega_{jj}}$ is continuously differentiable. Its partial derivatives are

$$\frac{\partial r_{ij}}{\partial \Omega_{ij}} = \frac{\text{sgn}(\Omega_{ij})}{\sqrt{\Omega_{ii}\Omega_{jj}}}, \quad \frac{\partial r_{ij}}{\partial \Omega_{ii}} = -\frac{|\Omega_{ij}|}{2\Omega_{ii}\sqrt{\Omega_{ii}\Omega_{jj}}},$$

and similarly for Ω_{jj} ; each is bounded in absolute value by C/η for a constant C depending on M . Thus the gradient norm $\|\nabla r_{ij}(\Omega)\|_1 \leq L$ for all $\Omega \in \mathcal{C}$. Since $\Omega^* \in \mathcal{C}$ and $\tilde{\Omega} \in \mathcal{C}$ whenever $\|\tilde{\Omega} - \Omega^*\|_\infty \leq 1$, the mean value theorem implies the claimed Lipschitz bound. \square

Proof of Lemma 5.4. Let $\delta_{ij} = \tilde{d}_{ij} - \tilde{c}_{ij}$ and $\delta_{ij}^* = d_{ij}^* - cr_{ij}^*$. From the triangle inequality and Lemma B.6,

$$|\delta_{ij} - \delta_{ij}^*| \leq |\tilde{d}_{ij} - d_{ij}^*| + c|\tilde{r}_{ij} - r_{ij}^*| \leq \varepsilon_B + cL\varepsilon_\Omega.$$

By bow-freeness, for each pair $\{i, j\}$ at most one of d_{ij}^* and r_{ij}^* is nonzero. Assumption 5.3 then implies $|\delta_{ij}^*| \geq \Delta^* > 0$.

Case 1: True directed edge ($d_{ij}^* > 0, r_{ij}^* = 0$). Then $\delta_{ij}^* = d_{ij}^* \geq \Delta^*$ and

$$\delta_{ij} \geq \delta_{ij}^* - (\varepsilon_B + cL\varepsilon_\Omega) > \frac{1}{2}\Delta^* > 0,$$

so the rule $\tilde{d}_{ij} \geq \tilde{c}_{ij}$ selects the directed channel. The threshold τ_B ensures that at least one of $|\tilde{B}_{ij}|, |\tilde{B}_{ji}|$ exceeds τ_B , so the directed edge is retained.

Case 2: True bidirected edge ($r_{ij}^* > 0, d_{ij}^* = 0$). Then $\delta_{ij}^* = -cr_{ij}^* \leq -\Delta^*$ and

$$\delta_{ij} \leq \delta_{ij}^* + (\varepsilon_B + cL\varepsilon_\Omega) < -\frac{1}{2}\Delta^* < 0,$$

so the rule selects the bidirected channel; τ_Ω ensures it survives thresholding.

Case 3: No edge ($d_{ij}^* = r_{ij}^* = 0$). The error bounds imply $|\tilde{B}_{ij}|, |\tilde{B}_{ji}| < \tau_B$ and $|\tilde{\Omega}_{ij}| < \tau_\Omega$, so no spurious edge is retained.

In all cases, the bow projection recovers the correct edge type (or absence), completing the proof. \square

C EXTENDED EXPERIMENTS

Real-world evaluation (Sachs Dataset) We evaluated DECOR on the Sachs protein signaling dataset (Sachs et al., 2005) ($n = 853, p = 11$), a standard benchmark for causal discovery. The ground truth is an 18-edge DAG. We compared against several baselines, including standard DAG learners (NOTEARS (Zheng et al., 2018), GOLEM (Ng et al., 2020), GES (Chickering, 2002a)), latent variable methods (DECAMF (Agrawal et al., 2023), DCD (?)), and LiNGAM (Shimizu et al., 2006). We also analyzed the effect of our post-hoc reconciliation step. We report results for **DECOR** (the output of the alternating optimization) and **DECOR (bow-strict)** (applying the reconciliation rule to enforce strict bow-freeness).

Table 1: Results on Sachs Data

Method	SHD	Precision	Recall	F1	Edges Pred.
DECOR	15	0.616	0.444	0.516	13
DECOR (bow-strict)	15	0.714	0.278	0.400	7
GES	16	0.571	0.444	0.500	14
DECOR_GL	18	0.500	0.167	0.250	6
DECOR_GL (bow-strict)	18	0.500	0.111	0.182	4
LiNGAM	19	0.455	0.222	0.348	11
DECAMF (LIN r1)	19	0.444	0.222	0.296	9
DCD	19	0.333	0.056	0.095	3
GOLEM	20	0.375	0.167	0.231	10
NOTEARS	22	0.167	0.056	0.083	6
AdaScore	25	0.333	0.389	0.359	21

Key Findings: DECOR achieves the lowest SHD (15) and highest F1 score (0.516), outperforming both classic methods (LiNGAM, GES) and recent latent-variable baselines (DCD, AdaScore). The `bow-strict` reconciliation acts as a rigorous filter: it significantly improves **Precision** (from 0.616 to 0.714, the highest among all methods) by enforcing the mutual exclusivity of directed and bidirected edges. However, this comes at the cost of **Recall** (dropping from 0.444 to 0.278) in this specific dataset, as weaker causal signals are conservatively assigned to the noise covariance. **DECOR vs. DECOR_GL:** The covariance-based formulation (DECOR) is more liberal, recovering more edges (13 predicted) compared to the precision-based formulation (DECOR_GL, 6 predicted), which appears over-penalized on this dataset.