

Supplementary Material for: Globally optimal score-based learning of directed acyclic graphs in high-dimensions

Appendix A contains some background and preliminary material that is important for the proof. We then use Appendix B to outline the main ideas and postpone detailed proofs of the various technical results to Appendices C-D. The reader interested in skipping directly to the proofs of the main theorems can find them in Appendices B.5 (Theorem 3.1) and B.6 (Theorem 4.1).

A Preliminaries

We begin by reviewing the connection between the equivalence class $\mathfrak{D}(\Sigma)$, Cholesky factors of Σ , and permutations. This material is essential and forms the basis for our proof technique. We then give some details on the regularizers used, and conclude with some important definitions on neighbourhood regression problems and introduce the concept of a *model selection exponent*.

A.1 Permutations

Denote the class of permutations on p elements by \mathbb{S}_p . For each $\pi \in \mathbb{S}_p$, define the associated permutation operator P_π on matrices: For any matrix A , $P_\pi A$ is the matrix obtained by permuting the rows and columns of A according to π , so that $(P_\pi A)_{ij} = a_{\pi(i)\pi(j)}$.

Cholesky representation Fix $\pi \in \mathbb{S}_p$. Write $\Gamma := \Sigma^{-1}$ and use the (modified) Cholesky decomposition to write $P_\pi \Gamma = (I - L)D^{-1}(I - L)^T$ where L is strictly lower triangular and $D \in \mathbb{R}_+^p$. Define $\tilde{B}(\pi) := P_{\pi^{-1}}L$ and $\tilde{\Omega}(\pi) := P_{\pi^{-1}}D$. The following result is well-known, but is restated here for completeness:

Lemma A.1. For any $\Sigma \succ 0$, the equivalence class of Σ (cf. (4)) is $\mathfrak{D}(\Sigma) = \{\tilde{B}(\pi) : \pi \in \mathbb{S}_p\}$.

Thus we can always write an arbitrary element of $\mathfrak{D}(\Sigma)$ as $\tilde{B}(\pi)$. The permutation π represents a valid topological sort for $\tilde{B}(\pi)$. The columns of $\tilde{B}(\pi)$ will be denoted by $\tilde{\beta}_j(\pi)$, and the j th diagonal element of $\tilde{\Omega}(\pi)$ will be denoted by $\tilde{\omega}_j^2(\pi)$. It follows from these definitions and (2) that

$$X_j = \tilde{\beta}_j(\pi)^T X + \tilde{\varepsilon}_j(\pi), \quad \text{where } \tilde{\varepsilon}_j(\pi) \sim \mathcal{N}(0, \tilde{\omega}_j^2(\pi)), \quad (15)$$

for $j = 1, \dots, p$. Note that $\text{supp}(\tilde{\beta}_j(\pi)) \subset S_j(\pi)$ for all $j = 1, \dots, p$, where

$$S_j(\pi) := \{k : \pi^{-1}(k) > \pi^{-1}(j)\} \quad (16)$$

consists of the nodes X_k that come after X_j under the ordering $X_{\pi(i)} \prec X_{\pi(i+1)}$ for $i = 1, \dots, p-1$.

Minimum-trace permutations In Section 2.1, we defined the notion of a minimum-trace DAG \tilde{B}_{\min} . By Lemma A.1, we know that $\tilde{B}_{\min} = \tilde{B}(\pi)$ for some $\pi \in \mathbb{S}_p$, where π is not necessarily unique. This motivates the following definition:

Definition A.1. The set of *minimum-trace permutations* is defined to be

$$\Pi_0 := \arg \min_{\pi \in \mathbb{S}_p} \text{tr} \tilde{\Omega}(\pi). \quad (17)$$

Given $\pi_0 \in \Pi_0$, π_0 is called a *minimum-trace permutation* and the corresponding DAG $\tilde{B}(\pi_0)$ is called a *minimum-trace DAG*. This definition does not require that \tilde{B}_{\min} is unique, and allows for the possibility that $\tilde{B}(\pi_1) \neq \tilde{B}(\pi_2)$ for $\pi_1, \pi_2 \in \Pi_0$.

Estimated permutations Recall that \mathbb{D} is the space of $p \times p$ real matrices that represent DAGs when interpreted as weighted adjacency matrices. For each $\pi \in \mathbb{S}_p$, define

$$\mathbb{D}[\pi] = \{B \in \mathbb{D} : P_\pi B \text{ is lower triangular}\}. \quad (18)$$

A DAG $B = [\beta_1 \mid \dots \mid \beta_p] \in \mathbb{D}$ is in $\mathbb{D}[\pi]$ if and only if $\text{supp}(\beta_j) \subset S_j(\pi)$ for all $j = 1, \dots, p$. In other words, for each node X_j , the permutation π defines a unique set of candidate parents given

538 by (16), and $B \in \mathbb{D}[\pi]$ if and only if the parent set of β_j comes from $S_j(\pi)$ for all j . By definition,
 539 $\tilde{B}(\pi) \in \mathbb{D}[\pi]$ for every π and hence $\text{supp}(\tilde{\beta}_j(\pi)) \subset S_j(\pi)$ for all j .

540 Recall the estimator \hat{B} defined via (1). The following definition formalizes the collection of permuta-
 541 tions that are topological sorts for \hat{B} :

542 *Definition A.2.* The collection of estimated permutations is

$$\hat{\Pi} = \{\pi \in \mathbb{S}_p : \hat{B} \in \mathbb{D}[\pi]\}.$$

543 An arbitrary element of $\hat{\Pi}$ will be denoted by $\hat{\pi}$. An equivalent definition of $\hat{\Pi}$ is the set of permutations
 544 π such that $P_\pi \hat{B}$ is lower triangular.

545 A.2 Regularizers

546 We study both the ℓ_1 and MCP regularizers as given by Condition 2.1. Here we summarize some
 547 properties of these regularizers for later use.

548 **Lemma A.2.** Suppose ρ_λ is either ℓ_1 or MCP. Then ρ_λ satisfies the following conditions:

549 (a) ρ_λ is concave and nondecreasing;

550 (b) $\rho_\lambda(0) = 0$;

551 (c) There are constants $\underline{\rho}_0, \underline{\rho}_1 \geq 0$, independent of λ , such that $\rho_\lambda(x) \geq \min\{\underline{\rho}_1 \lambda x, \underline{\rho}_0 \lambda^2\}$.

552 **Lemma A.3.** Suppose ρ_λ is either ℓ_1 or MCP. Then ρ_λ is additionally right-differentiable at zero
 553 and satisfies $0 < \rho'_\lambda(0+) < \infty$.

554 An elementary consequence of Lemma A.1 is that ρ_λ is subadditive. Lemma A.1(c) says that ρ_λ can
 555 be bounded below by a capped- ℓ_1 penalty: It is always true that a concave, nondecreasing function
 556 can be bounded below by a capped- ℓ_1 penalty, and Lemma A.1(c) simply normalizes this capped- ℓ_1
 557 penalty in terms of λ .

558 For completeness, we summarize below both regularizers under consideration along with the constants
 559 involved in the previous lemmas.

560 – The minimax concave penalty (MCP) proposed by Zhang [72]:

$$\rho_\lambda(x; \gamma) := \lambda \left(x - \frac{x^2}{2\lambda\gamma} \right) 1(x < \lambda\gamma) + \frac{\lambda^2\gamma}{2} 1(x \geq \lambda\gamma). \quad (19)$$

561 The MCP has $\rho'_\lambda(0+) = \lambda$, $\underline{\rho}_1 = 1/2$, and $\underline{\rho}_0 = \gamma/2$.

562 – The ℓ_1 penalty, $\rho_\lambda(x) = \lambda x$, has $\rho'_\lambda(0+) = \lambda$, $\underline{\rho}_1 = 1$, and $\underline{\rho}_0 \in [0, \infty)$.

563 Finally, since several of the results proved in this supplement do not require the incoherence condition
 564 $\zeta(G) < 1$, we will also make use of the following weaker version of Condition 2.1:

565 **Condition A.1** (Regularizer). The regularizer ρ_λ is chosen to be ℓ_1 or the MCP.

566 A.3 Neighbourhood regression

567 The core of our analysis is the regression decomposition (15) which we interpret as a neighbourhood
 568 regression problem and is used to learn the parent set of a node and hence the DAG structure. In this
 569 section we formalize these notions and introduce the concept of a *model selection exponent*, which
 570 quantifies the difficulty of a neighbourhood regression problem.

571 A.3.1 Penalized least-squares estimators

572 We are interested in the population SEM coefficients given by the following:

573 *Definition A.3.* For any $S \subset [p]_j$, let

$$\beta_j(S) := \arg \min_{\beta \in \mathbb{R}^p, \text{supp}(\beta) \subset S} \mathbb{E}[X_j - \beta^T X]^2.$$

574 We call $\beta_j(S)$ the SEM coefficients for X_j and denote the support set of $\beta_j(S)$ by $m_j(S) :=$
 575 $\text{supp}(\beta_j(S))$.

576 Note that $\beta_j(S) = \Sigma_{SS}^{-1} \Sigma_{Sj}$. Every positive definite matrix Σ defines a collection of $p2^{p-1}$ SEM
 577 coefficients given by $\{\beta_j(S) : S \subset [p]_j, j \in [p]\}$. We will be interested in estimating $\beta_j(S)$ via
 578 penalized least-squares (PLS):

579 *Definition A.4.* Suppose $y \in \mathbb{R}^n$ and $Z \in \mathbb{R}^{n \times m}$. Let $S \subset [m]$ and consider the set defined by

$$\widehat{\Theta}_\lambda(y, Z; S) := \arg \min_{\theta \in \mathbb{R}^m, \text{supp}(\theta) \subset S} \frac{1}{2n} \|y - Z\theta\|_2^2 + \rho_\lambda(\theta), \quad (20)$$

580 i.e., the set of global minimizers of the support-restricted PLS problem above. Let $\widehat{\Theta}_\lambda(y, Z) :=$
 581 $\widehat{\Theta}_\lambda(y, Z; [m])$ correspond to the case where there is no support restriction.

582 The support-restricted PLS problem $\widehat{\Theta}_\lambda(y, Z; S)$ allows us to properly define a neighbourhood
 583 regression problem. Let \mathbf{x}_j denote the j th column of \mathbf{X} .

584 *Definition A.5* (Neighbourhood regression). The *neighbourhood regression problem* for node X_j
 585 given a neighbourhood $S \subset [p]_j$ is defined to be the (possibly nonconvex) program given by
 586 $\widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S)$. An arbitrary solution to this program will be denoted by $\widehat{\beta}_j(S)$, i.e. $\widehat{\beta}_j(S) \in$
 587 $\widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S)$.

588 Learning \widehat{B} reduces to controlling $\widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S)$ for specific choices of S for each j (Lemma B.1).
 589 Thus, in the sequel, we no longer need to consider individual permutations, and instead will restrict
 590 our attention to subsets S , called candidate sets.

591 A.3.2 Model selection exponents

592 Given some $n \times m$ matrix Z and m -vector θ^* , define a set of “bad” noise vectors as follows:

$$A(Z, \theta^*; S) := \left\{ w \in \mathbb{R}^n : \text{supp}(\widehat{\theta}) \neq \text{supp}(\theta^*) \right. \\ \left. \text{for some } \widehat{\theta} \in \widehat{\Theta}_\lambda(Z\theta^* + w, Z; S) \right\}. \quad (21)$$

593 For a random vector $\mathbf{w} \in \mathbb{R}^n$ (e.g. $\mathbf{w} \sim \mathcal{N}_n(0, \sigma^2 I_n)$), we then have the following model selection
 594 failure event:

$$\mathcal{A}(\mathbf{w}, Z, \theta^*; S) := \{ \mathbf{w} \in A(Z, \theta^*; S) \}. \quad (22)$$

595 As usual we use the shorthand $\mathcal{A}(\mathbf{w}, Z, \theta^*) = \mathcal{A}(\mathbf{w}, Z, \theta^*; [m])$.

596 *Definition A.6.* Given a regularizer ρ_λ , the *model selection exponent* for the regression problem
 597 $\mathbf{y} = Z\theta^* + \mathbf{w}$ is defined to be

$$\Phi_\lambda(Z, \theta^*, \sigma^2) := -\log \mathbb{P} [\mathcal{A}(\mathbf{w}, Z, \theta^*)],$$

598 where \mathbb{P} is taken with respect to the distribution of $\mathbf{w} \sim \mathcal{N}_n(0, \sigma^2 I_n)$.

599 A larger exponent corresponds to better model selection performance. Let $\sigma_{\max}^2 :=$
 600 $\max_{1 \leq j \leq p} \text{var}(X_j)$ and note that $\sigma_{\max}^2 \leq r_{\max}(\Sigma)$. Define

$$\Psi_\lambda = \Psi_\lambda(\mathbf{X}) := \inf_{0 < \sigma \leq \sigma_{\max}} \inf_{\substack{\|\theta\|_0 \leq d(\Sigma) \\ \tau_*(\theta) \geq \tau_*(\Sigma)}} \Phi_\lambda(\mathbf{X}, \theta, \sigma^2). \quad (23)$$

601 This quantity measures “how difficult” the model selection problems defined by the fixed matrix \mathbf{X}
 602 are, and encodes what is usually proved in the regression literature: An upper bound on the probability
 603 of model selection failure given the maximum sparsity level d , minimum signal strength τ_* , and
 604 the maximum variance σ_{\max}^2 . This probability generally depends on the regularization parameter λ ,
 605 which in turn may depend on any of these quantities.

606 A.3.3 Example model selection exponents

607 To illustrate, let us derive a model selection exponent for the MCP, defined in (19). Huang et al. [27]
 608 consider PLS estimators $\hat{\Theta}_\lambda(y, Z; S)$ as defined in (20), applied to the data from a linear regression
 609 model $y = Z\theta^* + \mathbf{w}$, and provides conditions for model selection consistency. Adapting their result
 610 to our setup and notation, we have the following bound on model selection exponent for the MCP:

611 **Lemma A.4.** *Suppose $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$. Take $\rho_\lambda = \rho_\lambda(\cdot; \gamma)$ as in (19) and assume Σ is positive
 612 definite with bounded eigenvalues. Assume that*

- 613 1. $d(\Sigma) \leq \kappa_4 \cdot \min\{p, n, n/\log p\}$,
- 614 2. $\tau_*(\Sigma) > (1 + \gamma)\lambda$ for some $\gamma > \kappa_5 > 0$.

615 *Then for any $\lambda \geq \kappa_6 \cdot \sqrt{(d+1) \log p/n}$, it follows that $\mathbb{E}e^{-\Psi_\lambda(\mathbf{X}, \Sigma)} \leq 3 \exp(-2 \min\{d \log p, n\})$.
 616 Here, $\kappa_j = \kappa_j(\Sigma)$ ($j = 4, 5, 6$) are constants depending only on $\{r_{\min}(\Sigma), r_{\max}(\Sigma)\}$.*

617 This lemma is a straightforward consequence of Theorem 4.2 in [27] and Proposition 2 in [73].
 618 Briefly, [27] show that the least-squares MCP estimator correctly recovers the support of a linear
 619 model as long as the so-called *sparse Riesz condition* holds. We then use [73] to bound the probability
 620 that \mathbf{X} satisfies this condition. For the special case $\beta_j(S) = 0$ (which is not covered by [27]) we can
 621 invoke Proposition D.4.

622 In a similar manner, analogous bounds can be derived for other regularizers using existing results, see
 623 e.g. [27, 36, 66]. For example, using Corollary 1(a) in [36], a similar bound for ℓ_1 -regularization can
 624 be derived under the additional assumption that $\zeta(G) < 1$ as long as $n \gtrsim d \log p$.

625 B Outline of proofs

626 We seek control over the following event:

$$\mathcal{B} := \{\text{supp}(\hat{B}) \neq \text{supp}(\tilde{B}(\hat{\pi})) \exists \hat{\pi} \in \hat{\Pi}\}. \quad (24)$$

627 We will do this by reducing the analysis of \hat{B} to a family of neighbourhood regression problems.
 628 There are two key steps: (i) Showing that \hat{B} is equivalent to solving a series of p random regression
 629 problems given by $\hat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j(\hat{\pi}))$ (cf. Definition A.5), and (ii) Controlling the neighbourhood
 630 problems $\hat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j(\hat{\pi}))$ for all $\hat{\pi} \in \hat{\Pi}$.

631 The second step (ii) highlights the main technical difference between Theorems 3.1 and 4.1:

- 632 • To prove Theorem 3.1, we first prove a uniform concentration result for the score $Q(B)$, and
 633 use this to show that $\hat{\Pi} \subset \Pi_0$. That is, any estimated permutation must be a minimum-trace
 634 permutation. As a result, the random permutations $\hat{\pi} \in \hat{\Pi}$ are confined to live in a small set,
 635 which makes controlling the neighbourhood problems simpler. As a result we are able to
 636 bound $\mathbb{P}(\text{supp}(\hat{B}) \neq \text{supp}(\tilde{B}(\hat{\pi})))$ directly, which implies bounds on the desired quantities
 637 with $\hat{\pi}$ replaced by a minimum-trace permutation π_0 .
- 638 • To prove Theorem 4.1, we no longer assume we can restrict to a superstructure, and hence
 639 uniform score concentration (i.e. over the full space \mathbb{D}) is no longer readily viable. As a
 640 result, we must obtain *uniform control* over the neighbourhood problems $\hat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S)$ for
 641 all S and j . The challenge is that there are superexponentially many regression problems, so
 642 a naïve union bound over this family would yield overly pessimistic bounds on the order
 643 p/n . To deal with this, we will exploit a lattice property of these problems.

644 The proofs of Theorem 3.1 and 4.1 will be broken down into several steps. First, we establish
 645 some basic properties of the objective function and the probability space in order to reduce the
 646 neighbourhood regression analysis to a family of maximal sets denoted by $M_j(S)$ (Definition B.2).
 647 Then we introduce the lattice property (Lemmas B.2 and B.3) that is central to our proofs, and exploit
 648 this to provide a uniform bound on the probability of false selection for any neighbourhood problem
 649 (Proposition B.5).

650 **B.1 Reduction to neighbourhood regression**

651 Recall that the j th column of \widehat{B} is denoted by $\widehat{\beta}_j$ and denote the sample version of $\widetilde{\varepsilon}_j(\pi)$ by boldface,
 652 i.e. $\widetilde{\varepsilon}_j(\pi) := \mathbf{x}_j - \mathbf{X}\widehat{\beta}_j(\pi)$. The first step above is justified by the following result. The symbol $\perp\!\!\!\perp$
 653 is used here to denote independence of random variables.

654 **Lemma B.1.** *Suppose $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$ and $\lambda \geq 0$. Then the following statements are true:*

- 655 (a) *For any $j \in [p]$ and $\pi \in \mathbb{S}_p$, $\widetilde{\varepsilon}_j(\pi) \perp\!\!\!\perp \mathbf{X}_{S_j(\pi)}$.*
 656 (b) *\widehat{B} is a global minimizer of $Q(B)$ if and only if $\widehat{\beta}_j \in \widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j(\widehat{\pi}))$ for each $j \in [p]$ and
 657 $\widehat{\pi} \in \widehat{\Pi}$.*

658 The proof of this lemma, which is a simple consequence of how the least-squares loss and the
 659 regularizer factor, is found in Appendix C.3. This allows us to formally establish the equivalence
 660 between the DAG problem and neighbourhood regression: In order to construct \widehat{B} , it suffices to
 661 solve a neighbourhood regression problem for each column of \widehat{B} , given by $\widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j(\widehat{\pi}))$. A
 662 key observation is that through the independence established in Lemma B.1(a) and a conditioning
 663 argument, we can reduce the regression problem given by $\widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j(\widehat{\pi}))$ to a fixed design problem.
 664 The details are outlined in the proof of Proposition B.5.

665 **B.2 Invariant sets and monotonicity**

666 As a consequence of Lemma B.1, we have (cf. (24))

$$\mathcal{B} \subset \{\text{supp}(\widehat{\beta}_j(S)) \neq \text{supp}(\beta_j(S)) \exists j \in [p], S \subset [p]_j\}. \quad (25)$$

667 In order to further reduce the total number of estimators we must control, we will introduce the notion
 668 of an *invariant set*. First, recall the definition of $\beta_j(S)$ (cf. Definition A.3) and for any $j \in [p]$ and
 669 $S \subset [d]_j$ define the error (or noise) for the associated neighbourhood regression as the following
 670 residual:

$$\varepsilon_j(S) := X_j - \beta_j(S)^T X.$$

671 The support set of $\beta_j(S)$ is denoted by $m_j(S) := \text{supp}(\beta_j(S))$ and the error variance by $\omega_j^2(S) :=$
 672 $\text{var}(\varepsilon_j(S))$.

673 **Definition B.1.** For any $S \subset [p]_j$, define a collection of subsets by

$$\mathcal{T}_j(S) := \{T \subset [p]_j : \beta_j(T) = \beta_j(S)\} = \{T \subset [p]_j : m_j(T) = m_j(S)\},$$

674 where $\beta_j(S)$ and $m_j(S)$ are defined in Definition A.3. If $T \in \mathcal{T}_j(S)$, we call T an *invariant set of S*
 675 *for j* , or *S -invariant* for short.

676 In other words, for any j , $\mathcal{T}_j(S)$ is the collection of candidate sets $T \subset [p]_j$ such that the projection
 677 of X_j onto $\{X_i, i \in T\}$ is invariant. With some abuse of terminology, let us refer to $m_j(T) =$
 678 $\text{supp}(\beta_j(T))$ as the *support of neighbourhood T* (for node j). An equivalent description of $\mathcal{T}_j(S)$ is
 679 the set of neighbourhoods T whose support (for node j) is the same and equals $m_j(S)$.

680 The following lemma illustrates a crucial property of invariant sets:

681 **Lemma B.2.** $T_1, T_2 \in \mathcal{T}_j(S) \implies T_1 \cup T_2 \in \mathcal{T}_j(S)$.

682 This justifies the following definition:

683 **Definition B.2.** The unique largest element of $\mathcal{T}_j(S)$ shall be denoted by $M_j(S)$. Formally,

$$M_j(S) := \bigcup \mathcal{T}_j(S) = \bigcup \{T \subset [p]_j : \beta_j(T) = \beta_j(S)\}.$$

684 The name “ S -invariant set” comes from the fact that for any $T \in \mathcal{T}_j(S)$, we have the following useful
 685 identities:

$$\beta_j(m_j(S)) = \beta_j(S) = \beta_j(T) = \beta_j(M_j(S)), \quad (26)$$

$$\varepsilon_j(m_j(S)) = \varepsilon_j(S) = \varepsilon_j(T) = \varepsilon_j(M_j(S)). \quad (27)$$

686 The reason for introducing invariant sets is that it is generally *sufficient* to study the neighbourhood
687 problem for $M_j(S)$ in the sense that once we have model selection consistency for each estimator
688 in $\widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; M_j(S))$, the same is *guaranteed* for estimators based on every other neighbourhood
689 in $\mathcal{T}_j(S)$. In fact, we have the following result, which says that model selection properties of the
690 S -restricted estimators are monotone with respect to those sets S that contain the true support.

691 **Lemma B.3.** *Suppose that $Z \in \mathbb{R}^{n \times m}$ is fixed and consider the regression problem $\mathbf{y} = Z\theta^* + \mathbf{w}$
692 for some $\theta^* \in \mathbb{R}^m$. If $\text{supp}(\theta^*) \subset S \subset U$, then we have the following inclusion: $A(Z, \theta^*; S) \subset$
693 $A(Z, \theta^*; U)$. In particular, $\mathcal{A}(\mathbf{w}, Z, \theta^*; S) \subset \mathcal{A}(\mathbf{w}, Z, \theta^*; U)$ where $A(Z, \theta^*; S)$ and $\mathcal{A}(\mathbf{w}, Z, \theta^*; S)$
694 are defined in (21)–(22).*

695 We are interested in the model selection failure of $\widehat{\beta}_j(S)$ for $\beta_j(S)$, which can be stated as

$$\left\{ \text{supp}(\widehat{\beta}_j(S)) \neq \text{supp}(\beta_j(S)) \right. \\ \left. \text{for some } \widehat{\beta}_j(S) \in \widehat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S) \right\} = \mathcal{A}(\widetilde{\varepsilon}_j(S), \mathbf{X}, \beta_j(S); S) \quad (28)$$

696 in the notation introduced in (22).

697 **Corollary B.4.** *Suppose $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$. For any $S \subset [p]_j$, we have*

$$\mathcal{A}(\widetilde{\varepsilon}_j(S), \mathbf{X}, \beta_j(S); S) \subset \mathcal{A}(\widetilde{\varepsilon}_j(M_j(S)), \mathbf{X}, \beta_j(M_j(S)); M_j(S)).$$

698 Lemma B.4 is a *deterministic* statement about the events defined in (28), and proves that in order to
699 control the neighbourhood regression problem for some set $S \subset [p]_j$, it suffices to control the strictly
700 harder problem given by $M_j(S)$.

701 B.3 A bound on false selection

702 For any $\Sigma \succ 0$ and fixed node X_j , define the following collections of subsets:

$$m_j(\Sigma) := \{m_j(S) : S \subset [p]_j\}, \quad (29)$$

$$M_j(\Sigma) := \{M_j(S) : S \subset [p]_j\}. \quad (30)$$

703 Note that $|m_j(\Sigma)| = |M_j(\Sigma)|$. As long as it is clear whether the argument is a set S or a matrix Σ ,
704 this should not cause any confusion with $m_j(S)$ and $M_j(S)$.

705 For any neighbourhood $S \subset [p]_j$, recall that the associated error variance is given by $\omega_j^2(S) =$
706 $\text{var}(\varepsilon_j(S))$. With some more abuse of notation, let

$$\Phi_j(S) := \Phi_\lambda(\mathbf{X}_S, (\beta_j(S))_S, \omega_j^2(S)). \quad (31)$$

707 Note that we must restrict the SEM coefficients $\beta_j(S)$ to the subset S in order for this exponent to be
708 well-defined. Since $\text{supp}(\beta_j(S)) \subset S$, this does not change anything. The following general result
709 gives a uniform upper bound on the probability of false selection for any neighbourhood problem in
710 terms of the maximal sets $M_j(T)$.

711 **Proposition B.5.** *Fix $j \in [p]$ and $\Sigma \succ 0$. Then we have*

$$\mathbb{P}(\text{supp}(\widehat{\beta}_j(S)) \neq \text{supp}(\beta_j(S)), \exists S \subset [p]_j) \leq \sum_{T \in m_j(\Sigma)} \mathbb{E} e^{-\Phi_j(M_j(T))},$$

712 where $m_j(\Sigma)$ is defined by (29) and $\Phi_j(\cdot)$ is defined by (31).

713 The proof of this result can be found in Appendix C.7. The following result—which is proved in the
714 course of proving Proposition B.5—will also be useful when proving Theorem 3.1:

715 **Corollary B.6.** *Fix $j \in [p]$, $S \subset [p]_j$ and $\Sigma \succ 0$. Then we have*

$$\mathbb{P}(\text{supp}(\widehat{\beta}_j(T)) \neq \text{supp}(\beta_j(T)), \exists T \in \mathcal{T}_j(S)) \leq \mathbb{E} e^{-\Phi_j(M_j(S))},$$

716 where $\Phi_j(\cdot)$ is defined by (31).

717 Proposition B.5 says that to control the probability of false selection uniformly for all 2^{p-1} neigh-
718 bourhoods S of the node j , it suffices to control a much smaller class of problems given by the
719 neighbourhoods $M_j(T)$ for each support set $T \in m_j(\Sigma)$.

720 **B.4 Uniform support recovery**

721 The following result is a key ingredient in the proofs of both Theorem 3.1 and 4.1. It establishes an
722 upper bound on the probability of false selection, uniform over all S and j .

723 **Theorem B.1.** *Suppose $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$ with $\Sigma \succ 0$. Then*

$$\mathbb{P}(\text{supp}(\hat{\beta}_j(S)) = \text{supp}(\beta_j(S)), \forall j \in [p], S \subset [p]_j) \geq 1 - p \binom{p}{d} \mathbb{E} e^{-\Psi_\lambda(\mathbf{X}, \Sigma)}.$$

724 *Proof.* For any $T \in m_j(\Sigma)$, Lemma B.3 applied with $S = M_j(T)$ and $U = [p]$ yields

$$\Phi_j(M_j(T)) \geq \Phi_\lambda(\mathbf{X}, \beta_j(T), \omega_j^2(T)).$$

725 Recalling $d(\Sigma)$ and $\tau_*(\Sigma)$ in Definition 4.1, we have $\|\beta_j(T)\|_0 \leq d(\Sigma)$ and $\tau_*(\beta_j(T)) \geq \tau_*(\Sigma)$, as
726 well as $\omega_j^2(T) \leq \sigma_{\max}^2$. The previous expression combined with (23) implies:

$$\Phi_j(M_j(T)) \geq \Psi_\lambda(\mathbf{X}, \Sigma) \quad \text{for all } T \in M_j(\Sigma). \quad (32)$$

727 Combining Proposition B.5, (32) and a union bound over $j \in [p]$,

$$\begin{aligned} \mathbb{P}(\text{supp}(\hat{\beta}_j(S)) \neq \text{supp}(\beta_j(S)), \exists j \in [p], S \subset [p]_j) \\ \leq \sum_{j=1}^p \sum_{T \in m_j(\Sigma)} \mathbb{E} \exp(-\Phi_j(M_j(T))) \\ \leq p \binom{p}{d} \mathbb{E} \exp(-\Psi_\lambda(\mathbf{X}, \Sigma)), \end{aligned} \quad (33)$$

728 since there are at most $\binom{p}{d}$ subsets in $m_j(\Sigma)$. □

729 **B.5 Proof of Theorem 3.1**

730 Define $\ell(B) = \|\mathbf{X} - \mathbf{X}B\|_F^2 / (2n)$. There are two terms that we need to control: (i) The fluctuations
731 $|\ell(B) - \mathbb{E}\ell(B)|$ and (ii) The population loss $\mathbb{E}\ell(\hat{B})$. The fluctuations (i) are controlled by the
732 following proposition, which is proved in Appendix C.8, and may be of independent interest due to
733 its uniform control of an unbounded, subexponential empirical process:

734 **Proposition B.7.** *Let $\ell(B) = \|\mathbf{X} - \mathbf{X}B\|_F^2 / (2n)$ and let $\gamma_1(G)$ and $\gamma_2(G)$ be defined by (8) and
735 (9). Assume $\gamma_1 \leq 1$. Then there is a constant $\kappa(\Sigma; s)$, depending only on Σ and s , such that*

$$|\ell(B) - \mathbb{E}\ell(B)| \leq \gamma_1 [1 + 6\kappa(\Sigma; s)\gamma_2] \mathbb{E}\ell(B) \quad \text{for all } B \in \mathbb{D}_G \quad (34)$$

736 *with probability at least $1 - \alpha$, where*

$$\alpha := 2p^{-1} \left(\frac{9\gamma_1 p}{s} \right)^{-s} + \left(\frac{ep}{2s} \right)^{-s}. \quad (35)$$

737 For (ii), we have the following lemma:

738 **Lemma B.8.** *For any $\pi \in \mathbb{S}_p$, we have*

$$\mathbb{E}\ell(B) \geq \mathbb{E}\ell(\tilde{B}(\pi)) = \text{tr} \tilde{\Omega}(\pi) \quad \text{for all } B \in \mathbb{D}_p[\pi], \quad (36)$$

739 *where equality holds if and only if $B = \tilde{B}(\pi)$.*

740 Lemma B.8 implies, in particular, that $\mathbb{E}\ell(\hat{B}) \geq \mathbb{E}\ell(\tilde{B}(\hat{\pi}))$.

741 By Condition 3.1(a), we have $\tilde{B}(\pi_0) = \tilde{B}_{\min}$ and $\tilde{\Omega}(\pi_0) = \tilde{\Omega}_{\min}$ for all $\pi_0 \in \Pi_0$. For any two
742 permutations $\pi', \pi \in \mathbb{S}_p$ and $\eta > 0$, define a function

$$h(\pi', \pi; \eta) = (1 - \eta) \text{tr} \tilde{\Omega}(\pi') - (1 + \eta) \text{tr} \tilde{\Omega}(\pi) - \rho_\lambda(\tilde{B}(\pi)). \quad (37)$$

743 so that we have, recalling (7),

$$\chi(\eta) := \inf_{\pi' \notin \Pi_0} \sup_{\pi \in \Pi_0} h(\pi', \pi; \eta).$$

744 Recall that s is the maximum degree of G and define $\alpha = 2p^{-1}(9\gamma_1 p/s)^{-s} + (2p/s)^{-s}$. Since
 745 $\gamma_1 \leq 1$ (and hence $n \gtrsim s \log p$), Proposition B.7 and Lemma B.8 together imply that for any $\pi \in \Pi_0$
 746 and $\hat{\pi} \in \hat{\Pi}$, we have with probability $1 - \alpha$

$$\begin{aligned} (1 - \gamma_3) \mathbb{E} \ell(\tilde{B}(\hat{\pi})) + \rho_\lambda(\hat{B}) &\leq \ell(\hat{B}) + \rho_\lambda(\hat{B}) \\ &\leq \ell(\tilde{B}(\pi)) + \rho_\lambda(\tilde{B}(\pi)) \\ &\leq (1 + \gamma_3) \mathbb{E} \ell(\tilde{B}(\pi)) + \rho_\lambda(\tilde{B}(\pi)) \end{aligned}$$

747 where $\gamma_3 := \gamma_1[1 + 6\kappa(\Sigma; s)\gamma_2]$. Observing that $\mathbb{E} \ell(\tilde{B}(\pi)) = \text{tr } \tilde{\Omega}(\pi)$, we thus have $h(\hat{\pi}, \pi; \gamma_3) \leq 0$,
 748 where h is given by (37).

749 We now show that $\hat{\Pi} \subset \Pi_0$. Indeed, suppose $\hat{\pi} \notin \Pi_0$ for some $\hat{\pi} \in \hat{\Pi}$. Then by Condition 3.1(b),
 750 $\chi(\gamma_3) = \chi(\eta) > 0$, whence

$$\sup_{\pi \in \Pi_0} h(\hat{\pi}, \pi; \gamma_3) \geq \inf_{\pi' \notin \Pi_0} \sup_{\pi \in \Pi_0} h(\pi', \pi; \gamma_3) = \chi(\gamma_3) > 0,$$

751 which contradicts $h(\hat{\pi}, \pi; \gamma_3) \leq 0$. Thus $\hat{\Pi} \subset \Pi_0$.

752 By Lemma B.1(b) it suffices to show that

$$\mathbb{P}(\text{supp}(\hat{\beta}_j(S_j(\hat{\pi}))) \neq \text{supp}(\beta_j(S_j(\hat{\pi}))) \exists j \in [p]) = O(e^{-k \log p}). \quad (38)$$

753 Since the minimum-trace DAG is unique, for any $\pi, \pi' \in \Pi_0$, it follows that $m_j(S_j(\pi)) =$
 754 $\text{supp}(\tilde{\beta}_{\min, j}) = m_j(S_j(\pi'))$ and hence $M_j(S_j(\pi)) = M_j(S_j(\pi'))$. Using $\hat{\Pi} \subset \Pi_0$, we have

$$\begin{aligned} &\mathbb{P}(\text{supp}(\hat{\beta}_j(S_j(\hat{\pi}))) \neq \text{supp}(\beta_j(S_j(\hat{\pi}))) \exists j \in [p]) \\ &\leq \mathbb{P}(\text{supp}(\hat{\beta}_j(S_j(\pi_0))) \neq \text{supp}(\beta_j(S_j(\pi_0))) \exists j \in [p], \exists \pi_0 \in \Pi_0) \\ &= \mathbb{P}(\text{supp}(\hat{\beta}_j(S_j(\pi_0))) \neq \text{supp}(\tilde{\beta}_{\min, j}) \exists j \in [p], \exists \pi_0 \in \Pi_0) \\ &\leq \sum_{j=1}^p \mathbb{E} e^{-\Phi_j(M_j(\text{supp}(\tilde{\beta}_{\min, j})))}, \end{aligned}$$

755 where we used Corollary B.6 in the last line. Finally, apply known bounds (see Appendix A.3.3) to
 756 deduce $\Phi_j(M_j(\text{supp}(\tilde{\beta}_{\min, j}))) \gtrsim k \log p$ whenever $n \gtrsim k \log p$, which is implied since $k \leq s$ and
 757 we have assumed already that $n \gtrsim s \log p$. (If ℓ_1 -regularization is used, this is where we also need to
 758 assume $\zeta(G) < 1$ in Condition 2.1.) This implies the desired results with probability

$$1 - 2p^{-1} \left(\frac{9\gamma_1 p}{s} \right)^{-s} - \left(\frac{2p}{s} \right)^{-s} - O(e^{-k \log p}) = 1 - O(e^{-k \log p}),$$

759 where we used $k \leq s$ to simplify the probability bound. This completes the proof.

760 B.6 Proof of Theorem 4.1

761 The support recovery claim follows immediately from (25), Theorem B.1, and known bounds on
 762 the support recovery properties of penalized regression (see Section A.3.3 for discussion). Thus it
 763 remains to control $\rho_\lambda(\hat{B})$ and $\rho_\lambda(\tilde{B}(\hat{\pi}))$ by $\rho_\lambda(\tilde{B}(\pi_0))$.

764 The first step is the following lemma, which is a version of the standard basic inequality adapted to
 765 the current setting:

766 **Lemma B.9.** Let $\mathbf{E}(\pi) := \mathbf{X} - \mathbf{X}\tilde{B}(\pi)$. For any $\pi \in \mathbb{S}_p$ and $\hat{\pi} \in \hat{\Pi}$,

$$\begin{aligned} \frac{1}{2n} \|\mathbf{X}(\tilde{B}(\hat{\pi}) - \hat{B})\|_F^2 + \rho_\lambda(\hat{B}) &\leq \frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 - \frac{1}{2n} \|\mathbf{E}(\hat{\pi})\|_F^2 \\ &\quad + \frac{1}{n} \text{tr} \left(\mathbf{E}(\hat{\pi})^T \mathbf{X}(\tilde{B}(\hat{\pi}) - \hat{B}) \right) \\ &\quad + \rho_\lambda(\tilde{B}(\pi)). \end{aligned} \quad (39)$$

767 The proof of Lemma B.9 can be found in Appendix C.10. Lemma B.9 helps to reduce the analysis to
768 three terms:

769 **(B.9a)** The difference in residuals $\|\mathbf{E}(\pi)\|_F^2/(2n) - \|\mathbf{E}(\hat{\pi})\|_F^2/(2n)$ explains the origin of the
770 minimum-trace permutation: We would like to make $\|\mathbf{E}(\pi)\|_F^2/(2n)$ as small as possible
771 in order to minimize this difference. By standard concentration arguments, $\|\mathbf{E}(\pi)\|_F^2/n$
772 is close to its expectation, $\text{tr} \tilde{\Omega}(\pi)$. Hence, we choose π to minimize $\text{tr} \tilde{\Omega}(\pi)$. The details
773 of this argument are in Appendix D.3; the explicit upper bound we use is detailed in
774 Proposition D.8.

775 **(B.9b)** The quantity $\text{tr}(\mathbf{E}(\hat{\pi})^T \mathbf{X}(\tilde{B}(\hat{\pi}) - \hat{B}))/n$ can be bounded using the Gaussian width condition
776 (Definition D.1). There is a subtlety regarding whether to decompose this along rows or
777 columns; see Lemma D.6.

778 **(B.9c)** The penalty on \hat{B} can be replaced with $\rho_\lambda(\tilde{B}(\hat{\pi}))$ by showing that $\rho_\lambda(\hat{B}) \gtrsim \rho_\lambda(\tilde{B}(\hat{\pi}))$
779 (Lemma D.7).

780 Once we have establish control of these three terms (the details of which are found in Appendix D),
781 we can prove the following bound in terms of the constants δ (cf. Definition D.1) and a_2 (cf.
782 Condition 4.1):

783 **Proposition B.10.** Assume $n > 8(d+1)\log p$. Under Condition A.1 on ρ_λ , further assume

$$\tau_*(\mathfrak{D}(\Sigma)) \geq \tau_\lambda \left(\frac{2(1+\delta)}{1-3\delta} \right) \quad \text{for some } \delta \in (0, 1/3).$$

784 Let $\tilde{B}_{\min} = \tilde{B}(\pi_0)$ be a minimum-trace DAG satisfying Condition 4.1. Then

$$\frac{2\delta}{1-\delta} \rho_\lambda(\tilde{B}(\hat{\pi})) \stackrel{(i)}{\leq} \rho_\lambda(\hat{B}) \stackrel{(ii)}{\leq} \frac{2}{1-\delta} \left(1 + \frac{10}{a_2} \right) \rho_\lambda(\tilde{B}(\pi_0)), \quad (40)$$

785 with probability at least $1 - c_1 e^{-c_2 \min\{n, (d+1)\log p\}} - p \binom{p}{d} \mathbb{E} e^{-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta)}$.

786 The proof of Proposition B.10 follows from a series of standard concentration arguments (Ap-
787 pendix D), and can be found in Appendix D.4.

788 Finally, the desired bounds on $\rho_\lambda(\hat{B})$ and $\rho_\lambda(\tilde{B}(\hat{\pi}))$ follow from Proposition B.10 by taking $\delta =$
789 $(a_1 - 2)/(3a_1 + 2) \in (0, 1/3)$, and using Proposition D.3 to complete the probability bound.

790 C Proofs of technical results

791 C.1 Proof of Lemma 2.1

792 Consider the following program:

$$\min \sum_{j=1}^p x_j^2 \quad \text{subject to} \quad \sum_{j=1}^p \log x_j^2 = C. \quad (41)$$

793 The solution to this program is given by $x_j^2 = e^{C/p}$ for all $j = 1, \dots, p$. In other words, the minimum
794 is attained by a constant vector. It is straightforward to verify that $\log \det \tilde{\Omega}(\pi) = \log \det \Sigma$ and
795 hence $\log \det \tilde{\Omega}(\pi) = \sum_j \log \tilde{\omega}_j^2(\pi)$ is constant for all $\pi \in \mathbb{S}_p$. Thus for any $\pi \in \mathbb{S}_p$, the vector
796 $(\tilde{\omega}_1^2(\pi), \dots, \tilde{\omega}_p^2(\pi)) \in \mathbb{R}^p$ is feasible for (41), which implies that $\text{tr} \tilde{\Omega}(\pi)$ is minimized whenever
797 $\tilde{\omega}_1^2(\pi) = \dots = \tilde{\omega}_p^2(\pi)$. Finally, uniqueness of $\tilde{B}(\pi_0)$ follows from Theorem 1 in Peters and Bühlmann
798 [47].

799 **C.2 Proof of Lemma A.1**

800 We need the following simple lemma, which follows since $P_\pi A = PAP^T$ for some permutation
801 matrix P :

802 **Lemma C.1.** $A = MNM^T \iff P_\pi A = (P_\pi M)(P_\pi N)(P_\pi M)^T$.

803 Recall the modified Cholesky decomposition of A (also called the LDLT decomposition): $A = LDL^T$
804 for a lower triangular matrix L , with unit diagonal entries, and a diagonal matrix D . When A is
805 positive definite, the pair (L, D) is unique and we refer to it as the *Cholesky decomposition of A* .

806 Let us denote the set of all pairs $(\tilde{B}, \tilde{\Omega})$ satisfying $\Sigma^{-1} = (I - \tilde{B})\tilde{\Omega}^{-1}(I - \tilde{B}^T)$ (equivalently, (3)) as
807 \mathfrak{D}' . Next, note that $\tilde{B} \in \mathbb{D}$ if and only if $P_\pi \tilde{B}$ is lower triangular for some permutation π . Lemma C.1
808 implies that $(\tilde{B}, \tilde{\Omega}) \in \mathfrak{D}'$ iff $(I - P_\pi \tilde{B}, P_\pi \tilde{\Omega}^{-1})$ is a Cholesky decomposition of $P_\pi \Sigma^{-1}$ for some π .

809 Now, $(I - P_\pi \tilde{B}(\pi), P_\pi \tilde{\Omega}(\pi)^{-1})$ is also a Cholesky decomposition of $P_\pi \Sigma^{-1}$. Since the Cholesky de-
810 composition is unique for positive definite matrices, we have $(\tilde{B}, \tilde{\Omega}) \in \mathfrak{D}'$ iff $(\tilde{B}, \tilde{\Omega}) = (\tilde{B}(\pi), \tilde{\Omega}(\pi))$
811 for some π , which gives the desired result, since $\mathfrak{D}(\Sigma)$ is the projection of \mathfrak{D}' onto its first coordinate.

812 **C.3 Proof of Lemma B.1**

813 The first conclusion (a) follows from elementary properties of conditional expectation and the identity

$$\mathbb{E}(X_j | X_{S_j(\pi)}) = \tilde{\beta}_j(\pi)^T X.$$

814 To prove (b), fix $\hat{\pi} \in \hat{\Pi}$ and let $S_j = S_j(\hat{\pi})$. If $\hat{\beta}_j \in \hat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j)$ for each j , then evidently
815 $\hat{B} = [\hat{\beta}_1 | \dots | \hat{\beta}_p]$ minimizes $Q(B)$ over $\mathbb{D}[\hat{\pi}]$ (cf. (18)). For the reverse direction, recall that
816 \mathbf{X}_{S_j} is the $n \times |S_j|$ matrix formed by extracting the columns in S_j , and similarly for $(\beta_j)_{S_j}$. For
817 any $B \in \mathbb{D}[\pi]$ we have $(\beta_j)_{S_j^c} = 0$ for each j , so we can write fix $\hat{\pi} \in \hat{\Pi}$ and let $S_j = S_j(\hat{\pi})$. If
818 $\hat{\beta}_j \in \hat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j)$ for each j , then evidently $\hat{B} = [\hat{\beta}_1 | \dots | \hat{\beta}_p]$ minimizes $Q(B)$ over $\mathbb{D}[\hat{\pi}]$. For
819 the reverse direction, recall that \mathbf{X}_{S_j} is the $n \times |S_j|$ matrix formed by extracting the columns in S_j ,
820 and similarly for $(\beta_j)_{S_j}$. For any $B \in \mathbb{D}[\pi]$ we have $(\beta_j)_{S_j^c} = 0$ for each j , so we can write

$$\begin{aligned} \frac{1}{2n} \|\mathbf{X} - \mathbf{X}B\|_F^2 + \rho_\lambda(B) &= \sum_{j=1}^p \left\{ \frac{1}{2n} \|\mathbf{x}_j - \mathbf{X}\beta_j\|_2^2 + \rho_\lambda(\beta_j) \right\} \\ &= \sum_{j=1}^p \left\{ \frac{1}{2n} \|\mathbf{x}_j - \mathbf{X}_{S_j}(\beta_j)_{S_j}\|_2^2 + \rho_\lambda((\beta_j)_{S_j}) \right\}. \end{aligned}$$

821 Then $\hat{B} \in \min_{\mathbb{D}[\hat{\pi}]} Q(B)$ if and only if

$$\hat{\beta}_j \in \arg \min_{\beta} \frac{1}{2n} \|\mathbf{x}_j - \mathbf{X}\beta\|_2^2 + \rho_\lambda(\beta) \quad \text{subject to } \beta_{S_j^c} = 0.$$

822 In other words, $\hat{\beta}_j \in \hat{\Theta}_\lambda(\mathbf{x}_j, \mathbf{X}; S_j)$ for each j . Since $\hat{\pi} \in \hat{\Pi}$ was arbitrary, the desired claim follows.

823 **C.4 Proof of Lemma B.2**

824 The proof relies on the following property of L^2 projections: For any two sets $S, R \subset [p]_j$, we have

$$\beta_j(S \cup R) = \beta_j(S) \iff \varepsilon_j(S) \perp\!\!\!\perp X_i, \forall i \in R. \quad (42)$$

825 To lighten the notation, let $S^* = m_j(S)$. Note that $\beta_j(S) = \beta_j(S^*)$ since $\text{supp}(\beta_j(S)) = S^*$. It
826 follows from (42) that $\varepsilon_j(S^*) \perp\!\!\!\perp X_i$ for $i \in S \setminus S^*$. Similarly, since $\text{supp}(\beta_j(T_k)) = S^*$, we have
827 $\varepsilon_j(S^*) \perp\!\!\!\perp X_i$ for $i \in T_k \setminus S^*$ and $k = 1, 2$. It follows that

$$\varepsilon_j(S^*) \perp\!\!\!\perp X_i, \forall i \in (T_1 \setminus S^*) \cup (T_2 \setminus S^*)$$

828 hence the application of (42) in the reverse direction yields

$$\beta_j(T_1 \cup T_2) = \beta_j(S^* \cup (T_1 \setminus S^*) \cup (T_2 \setminus S^*)) = \beta_j(S^*) = \beta_j(S).$$

829 **C.5 Proof of Lemma B.3**

It suffices to show

$$A(Z, \theta^*; U)^c \subset A(Z, \theta^*; S)^c.$$

830 Suppose $w \in A(Z, \theta^*; U)^c$, i.e., $\text{supp}(\tilde{\theta}) = \text{supp}(\theta^*) := S^*$ for any $\tilde{\theta} \in \hat{\Theta}_\lambda(Z\theta^* + w, Z; U)$. We
831 wish to show that for any $\hat{\theta} \in \hat{\Theta}_\lambda(Z\theta^* + w, Z; S)$, it must also be true that $\text{supp}(\hat{\theta}) = S^*$. Let

$$F(\theta) = \frac{1}{2n} \|Z(\theta^* - \theta) + w\|_2^2 + \rho_\lambda(\theta)$$

832 denote the objective function in Definition A.4 of $\hat{\Theta}_\lambda(y, Z; S)$ with $y = Z\theta^* + w$. Since $\text{supp}(\hat{\theta}) \subset$
833 $S \subset U$, $\hat{\theta}$ is feasible for the U -restricted problem, whence

$$F(\tilde{\theta}) \leq F(\hat{\theta})$$

834 for any $\tilde{\theta} \in \hat{\Theta}_\lambda(Z\theta^* + w, Z; U)$. But $\tilde{\theta}$ is also feasible for the S -restricted problem since $\text{supp}(\tilde{\theta}) =$
835 $S^* \subset S$, so that

$$F(\tilde{\theta}) \geq F(\hat{\theta}) \implies F(\tilde{\theta}) = F(\hat{\theta}).$$

836 Since the value $F(\tilde{\theta})$ is by definition the global minimum of F for the U -restricted problem and
837 $\text{supp}(\hat{\theta}) \subset U$, $\hat{\theta}$ must be a global minimizer of F for the U -restricted problem, i.e., $\hat{\theta} \in \hat{\Theta}_\lambda(Z\theta^* +$
838 $w, Z; U)$, whence $\text{supp}(\hat{\theta}) = S^*$ as desired.

839 **C.6 Proof of Corollary B.4**

840 By Lemma B.3 and the fact that $S \subset M_j(S)$, we have

$$\mathcal{A}(\tilde{\varepsilon}_j(S), \mathbf{X}, \beta_j(S); S) \subset \mathcal{A}(\tilde{\varepsilon}_j(S), \mathbf{X}, \beta_j(S); M_j(S)). \quad (43)$$

841 Using (26) and (27), we have the following identity:

$$\mathcal{A}(\tilde{\varepsilon}_j(S), \mathbf{X}, \beta_j(S); M_j(S)) = \mathcal{A}(\tilde{\varepsilon}_j(M_j(S)), \mathbf{X}, \beta_j(M_j(S)); M_j(S)).$$

842 Plugging this into (43) yields the desired result.

843 **C.7 Proof of Proposition B.5**

844 Throughout, for simplicity, let

$$\mathcal{A}_S := \mathcal{A}(\tilde{\varepsilon}_j(S), \mathbf{X}, \beta_j(S); S).$$

845 Fix $S \subset [p]_j$ and let $\theta^* = \beta_j(S)$, $s^* = |m_j(S)| = \|\theta^*\|_0$ and $\varepsilon^* = \tilde{\varepsilon}_j(S)$ so that $\mathcal{A}_S =$
846 $\mathcal{A}(\varepsilon^*, \mathbf{X}, \theta^*; S)$. Note that $\mathcal{A}(\varepsilon^*, \mathbf{X}, \theta^*; S)$ represents the following model selection failure:

$$\text{supp}(\hat{\theta}) \neq \text{supp}(\theta^*) \quad \exists \hat{\theta} \in \hat{\Theta}_\lambda(\mathbf{X}\theta^* + \varepsilon^*, \mathbf{X}; S).$$

847 Since $\text{supp}(\theta^*) \subset S$, we can restrict \mathbf{X} and θ^* to S , so that the above is equivalent to

$$\text{supp}(\hat{\theta}) \neq \text{supp}(\theta_S^*) \quad \exists \hat{\theta} \in \hat{\Theta}_\lambda(\mathbf{X}_S\theta_S^* + \varepsilon^*, \mathbf{X}_S).$$

848 which is the same event as $\mathcal{A}(\varepsilon^*, \mathbf{X}_S, \theta_S^*)$. To summarize, $\mathcal{A}_S = \mathcal{A}(\varepsilon^*, \mathbf{X}_S, \theta_S^*)$.

849 Since ε^* is independent of \mathbf{X}_S by Lemma B.1(a), by conditioning on \mathbf{X}_S we are dealing with a fixed
850 design regression problem with Gaussian noise $\varepsilon^* = \tilde{\varepsilon}_j(S) \sim \mathcal{N}_n(0, \omega_j^2(S)I_n)$. We obtain

$$\begin{aligned} \mathbb{P}(\mathcal{A}_S) &= \mathbb{E} \left[\mathbb{P}(\mathcal{A}(\varepsilon^*, \mathbf{X}_S, \theta_S^*) \mid \mathbf{X}_S) \right] \\ &\leq \mathbb{E} \exp[-\Phi_\lambda(\mathbf{X}_S, \theta_S^*, \omega_j^2(S))] \\ &= \mathbb{E} \exp(-\Phi_j(S)), \end{aligned} \quad (44)$$

851 where the last line uses (31). Now we have

$$\{\text{supp}(\widehat{\beta}_j(T)) \neq \text{supp}(\beta_j(T)), \exists T \in \mathcal{T}_j(S)\} = \bigcup_{T \in \mathcal{T}_j(S)} \mathcal{A}_T = \mathcal{A}_{M_j(S)}, \quad (45)$$

852 where the first equality is by (28) and the second follows from Corollary B.4. Note that this is the key
853 step where the reduction occurs. Hence, combining (45) with (44) we have

$$\begin{aligned} \mathbb{P}\left(\bigcup_{S \subset [p]_j} \mathcal{A}_S\right) &= \mathbb{P}\left(\bigcup_{S \subset [p]_j} \mathcal{A}_{M_j(S)}\right) \\ &= \mathbb{P}\left(\bigcup_{T \in m_j(\Sigma)} \mathcal{A}_{M_j(T)}\right) \\ &\leq \sum_{T \in m_j(\Sigma)} \mathbb{P}(\mathcal{A}_{M_j(T)}) \leq \sum_{T \in m_j(\Sigma)} \mathbb{E} \exp(-\Phi_j(M_j(T))), \end{aligned}$$

854 which is the desired probability bound.

855 C.8 Proof of Proposition B.7

856 We work with the column decomposition of the loss

$$\begin{aligned} \ell_j(\beta) &:= \frac{1}{n} \|\mathbf{X}(e_j - \beta)\|_2^2, \\ \mathbb{E} \ell_j(\beta) &= \frac{1}{2} (e_j - \beta)^T \Sigma (e_j - \beta), \end{aligned}$$

857 where $e_j \in \mathbb{R}^p$ is the j th standard basis vector. The overall loss can be written as

$$\ell(B) = \sum_{j=1}^n \ell_j(\beta_j)$$

858 where β_j is the j th column of B . Let us also define so that

$$J(\beta) = \frac{1}{n} \|\mathbf{X}\beta\|_2^2, \quad \text{and,} \quad \mathbb{E} J(\beta) := \beta^T \Sigma \beta$$

859 so that $\ell_j(\beta) = \frac{1}{2} J(e_j - \beta)$ and $\mathbb{E} \ell_j(\beta) = \frac{1}{2} \mathbb{E} J(e_j - \beta)$. It is easier to work with J . Let \mathbf{X}_i^T be the
860 i th row of \mathbf{X} . Then, $J(\beta) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{X}_i^T \beta)^2$.

861 Let $K = \mathbb{D}_G$ and K_j denote the set of β_j for $B \in K$. Define

$$\mathbb{B}_0^{-j}(s) = \{x \in \mathbb{R}^p : \|x\| \leq s, x_j = 0\}. \quad (46)$$

862 Note that $\beta_j \in \mathbb{B}_0^{-j}(s)$ for every $B \in \mathbb{D}_G$ and in particular $\widetilde{\beta}_j(\pi) \in \mathbb{B}_0^{-j}(s)$ if $\widetilde{B}(\pi) \in \mathbb{D}_G$. Finally,
863 define

$$\kappa(\Sigma; s) := \frac{\|\Sigma\|_{(2s+2)}}{r_{\min}^{(s+1)}(\Sigma)} \quad (47)$$

864 where

$$\|\Sigma\|_{(s)} := \max_{S: |S|=s} \|\Sigma_S\|, \quad (48)$$

$$r_{\min}^{(s)}(\Sigma) := \inf_{S: |S|=s} r_{\min}(\Sigma_S). \quad (49)$$

865 *Proof of Proposition B.7.* For any $\beta \in K_j$, we have

$$\|e_j - \beta\|^2 = 1 + \|\beta\|^2. \quad (50)$$

866 For any $t \leq 1$, applying Lemma C.5 we have that on the event \mathcal{B}_{2s} (defined in (53))

$$|\ell_j(\beta) - \mathbb{E} \ell_j(\beta)| \leq t J(\beta) + 3\gamma_2 (1 + \|\beta\|_2^2) \varepsilon \|\Sigma\|_{(2s)} \text{ for all } \beta \in K_j \quad (51)$$

867 fails with probability at most $2\left(\frac{3ep}{s\varepsilon}\right)^s e^{-cnt^2}$. A further union bound over $j = 1, \dots, p$ gives that

$$|\ell_j(\beta) - \mathbb{E}\ell_j(\beta)| \leq tJ(\beta) + 3\gamma_2(1 + \|\beta\|_2^2)\varepsilon\|\Sigma\|_{(2s)}, \quad (52)$$

for all $\beta \in K_j$ and all $j \in [p]$

868 fails with probability at most

$$2p\left(\frac{3ep}{s\varepsilon}\right)^s e^{-cnt^2} + \mathbb{P}(\mathcal{B}_{2s}^c) \leq 2p\left(\frac{3ep}{s\varepsilon}\right)^s e^{-cnt^2} + \left(\frac{ep}{2s}\right)^{-c_1 2s} =: T_1 + T_2,$$

869 where we invoked Lemma C.4 to bound $\mathbb{P}(\mathcal{B}_{2s}^c)$. Take $t = \varepsilon$ and let $N = p^{1/s}3ep/s$ so that the first
870 term in the bound is

$$T_1 := 2\left(\frac{N}{\varepsilon}\right)^s e^{-cn\varepsilon^2}.$$

871 Take ε

$$\varepsilon^2 = \frac{2}{c} \frac{s}{n} \log N \leq 1$$

872 which gives the following bound,

$$T_1 \leq 2(\varepsilon N)^{-s} = 2p^{-1}\left(\gamma_1 \frac{3ep}{s}\right)^{-s}.$$

873 where we note that with our choices, we have $t = \varepsilon = \gamma_1$ as defined in (8).

874 Note that for any $\beta \in K_j$, $(1 + \|\beta\|^2)r_{\min}^{(s+1)}(\Sigma) \leq 2\ell_j(\beta)$, which implies $1 + \|\beta\|^2 \leq$
875 $2\ell_j(\beta)/r_{\min}^{(s+1)}(\Sigma)$. Plugging this upper bound into (52) and summing over j gives (34), where
876 $\kappa(\Sigma; s)$ is defined by (47). The proof is complete. \square

877 Below we prove the various technical lemmas required in the previous proof.

878 **Lemma C.2.** *We have*

$$\mathbb{P}\left(|J(\beta) - \mathbb{E}J(\beta)| \geq t \cdot J(\beta)\right) \leq 2 \exp\left[-\frac{1}{8}n \cdot \min(t^2, t)\right], \quad t \geq 0.$$

879 *Proof.* Note that $\mathbf{X}_i^T \beta / \sqrt{\mathbb{E}J(\beta)} \sim N(0, 1)$ iid for each $i = 1, \dots, n$. Then,

$$\frac{J(\beta)}{\mathbb{E}J(\beta)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{\mathbf{X}_i^T \beta}{\sqrt{\mathbb{E}J(\beta)}}\right)^2 \sim \chi_n^2$$

880 and the claim follows from χ^2 concentration. \square

881 Let \mathcal{B}_{2s} be the following event:

$$\mathcal{B}_{2s} := \left\{g_{n,s}(\varepsilon) \leq \varepsilon \sqrt{\gamma_2} \|\Sigma\|_{(2s)}^{1/2}, \quad \forall \varepsilon > 0\right\}. \quad (53)$$

882 where

$$g_{n,s}(\varepsilon) = \sup_{\|u\|_2 \leq \varepsilon, \|u\|_0 \leq s} \frac{1}{\sqrt{n}} \|\mathbf{X}u\|_2, \quad \text{and} \quad (54)$$

$$\|\Sigma\|_{(s)} := \max_{S:|S|=s} \|\Sigma_S\|. \quad (55)$$

883 **Lemma C.3.** *For any $c_1 > 0$, with probability at least $1 - (ep/s)^{-c_1 s} e^{-u^2/2}$,*

$$g_{n,s}(\varepsilon) \leq \varepsilon \|\Sigma\|_{(s)}^{1/2} \left(1 + C \sqrt{\frac{s \log(ep/s)}{n}} + \frac{u}{\sqrt{n}}\right), \quad \forall \varepsilon > 0.$$

884 where $C = \sqrt{2(c_1 + 1)} + 1$.

885 *Proof.* Note that $g_{n,s}(\varepsilon) = \varepsilon g_{n,s}(1)$ which is obtained by the change of variable $u \rightarrow \varepsilon u$. We have

$$g_{n,s}(1) = \max_{S:|S|=s} \sup_{\|u\|_2 \leq 1} \frac{1}{\sqrt{n}} \|\mathbf{X}_S u\|_2.$$

886 Let $W_S = \mathbf{X}_S \Sigma_S^{-1/2}$ so that $W_S \sim N(0, I_s)$. Then

$$\mathbb{P}(\|W_S\| > \sqrt{s} + \sqrt{n} + t) \leq \exp(-t^2/2).$$

887 We also have $\|\mathbf{X}_S u\|_2 = \|\Sigma_S^{1/2} W_S u\|_2 = \|\Sigma_S^{1/2}\| \|W_S\| \|u\|_2$. Thus,

$$g_{n,s}(1) \leq \max_{S:|S|=s} \|\Sigma_S^{1/2}\| \cdot \max_{S:|S|=s} \frac{1}{\sqrt{n}} \|W_S\|$$

888 and hence

$$\mathbb{P}\left(\max_{S:|S|=s} \frac{1}{\sqrt{n}} \|W_S\| > \sqrt{\frac{s}{n}} + 1 + \frac{t}{\sqrt{n}}\right) \leq \binom{p}{s} \exp(-t^2/2)$$

889 Taking $t = \sqrt{2(c_1 + 1)s \log(ep/s)} + u$, the above probability is bounded by

$$\binom{p}{s} e^{-t^2/2} \leq (ep/s)^s (ep/s)^{-(c_1+1)s} e^{-u^2/2} = (ep/s)^{-c_1 s} e^{-u^2/2}$$

890 Letting $C = \sqrt{2(c_1 + 1)} + 1$ and noting that $\|A^{1/2}\| = \|A\|^{1/2}$, the result follows. \square

891 **Lemma C.4.** Let \mathcal{B}_{2s} be defined as in (53). Then $\mathbb{P}(\mathcal{B}_{2s}^c) \leq (ep/(2s))^{-c_1 2s}$.

892 *Proof.* Apply Lemma C.3 with s replaced with $2s$ and $u = 0$, and note that

$$2s \log(ep/(2s)) = 2s[\log(ep/s) - \log 2] \leq 2s \log(ep/s),$$

893 we observe that \mathcal{B}_{2s} fails with probability at most $(ep/(2s))^{-c_1 2s}$. \square

894 Define the L^q balls

$$\mathbb{B}_q(r) := \{x \in \mathbb{R}^p : \|x\|_q \leq r\}.$$

895 **Lemma C.5.** For any $t \leq 1$: On the event \mathcal{B}_{2s} ,

$$|J(\beta) - \mathbb{E}J(\beta)| \leq t J(\beta) + 3\gamma_2 \|\beta\|_2^2 \varepsilon \|\Sigma\|_{(2s)}, \quad \text{for all } \beta \in \mathbb{B}_0(s) \quad (56)$$

896 fails with probability at most $2\left(\frac{3ep}{s\varepsilon}\right)^s e^{-cnt^2}$.

897 *Proof.* Write (56) in the form $g(\beta) \leq 0$ and observe that g is homogeneous of order two: $g(r\beta) = r^2 g(\beta)$ for any $r \in \mathbb{R}$ and $\beta \in \mathbb{R}^p$. Thus, it is enough to establish the bound for $\beta \in \mathbb{B}_2(1)$. The general case is then obtained by applying the bound to $\beta/\|\beta\|_2$.

900 Let $N_j \subseteq \mathbb{B}_0(s) \cap \mathbb{B}_2(1)$ be an ε -net for $\mathbb{B}_0(s) \cap \mathbb{B}_2(1)$ in ℓ_2 norm. Then, $|N_j| \leq \binom{p}{s} (3/\varepsilon)^s$. For any $\beta, \beta' \in \mathbb{B}_0(s) \cap \mathbb{B}_2(1)$, using

$$\left| \|x\|^2 - \|y\|^2 \right| \leq \|x\| - \|y\| (\|x\| + \|y\|) \leq \|x - y\| \|x + y\|$$

902 we have on event \mathcal{B}_{2s} ,

$$\begin{aligned} |J(\beta) - J(\beta')| &= \left| \frac{1}{n} \|\mathbf{X}\beta\|_2^2 - \frac{1}{n} \|\mathbf{X}\beta'\|_2^2 \right| \\ &\leq \frac{1}{n} \|\mathbf{X}(\beta - \beta')\|_2 \cdot \|\mathbf{X}(\beta + \beta')\|_2 \\ &\leq g_{n,2s}(\|\beta - \beta'\|_2) \cdot g_{n,2s}(\|\beta + \beta'\|_2) \\ &\leq \sqrt{\gamma_2} \|\Sigma\|_{(2s)}^{1/2} \|\beta - \beta'\|_2 \cdot \sqrt{\gamma_2} \|\Sigma\|_{(2s)}^{1/2} \|\beta + \beta'\|_2 \\ &\leq 2\gamma_2 \|\Sigma\|_{(2s)} \|\beta - \beta'\|_2, \end{aligned}$$

903 where we have used $\|\beta + \beta'\|_2 \leq \|\beta\| + \|\beta'\| \leq 2$. A similar bound holds for the expectation $\mathbb{E}J(\beta)$,
 904 i.e.

$$\begin{aligned} |\mathbb{E}J_j(\beta) - \mathbb{E}J_j(\beta')| &= \left| \|\Sigma^{1/2}\beta\|_2^2 - \|\Sigma^{1/2}\beta'\|_2^2 \right| \\ &\leq \|\Sigma^{1/2}(\beta - \beta')\|_2 \cdot \|\Sigma^{1/2}(\beta + \beta')\|_2 \\ &\leq \|\Sigma\|_{(2s)}^{1/2} \|\beta - \beta'\|_2 \cdot \|\Sigma\|_{(2s)}^{1/2} \|\beta + \beta'\|_2 \\ &\leq 2\|\Sigma\|_{(2s)} \|\beta - \beta'\|_2. \end{aligned}$$

905 It follows that for any $f(\cdot)$, with $\gamma = 2(\gamma_2 + 1)$,

$$\begin{aligned} \sup_{\beta \in \mathbb{B}_0(s) \cap \mathbb{B}_2(1)} (|J(\beta) - \mathbb{E}J(\beta)| - f(\beta)) \\ \leq \gamma \|\Sigma\|_{(2s)} \varepsilon + \sup_{\beta \in N_j} (|J(\beta) - \mathbb{E}J(\beta)| - f(\beta)). \end{aligned} \quad (57)$$

906 We now have, for any $t \geq 0$,

$$\begin{aligned} \mathbb{P} \left(\left\{ \sup_{\beta \in \mathbb{B}_0(s) \cap \mathbb{B}_2(1)} \left[|J(\beta) - \mathbb{E}J(\beta)| - t\mathbb{E}J(\beta) \right] \geq \gamma \|\Sigma\|_{(2s)} \varepsilon \right\} \cap \mathcal{B}_{2s} \right) \\ \leq \mathbb{P} \left(\left\{ \sup_{\beta \in N_j} \left[|J(\beta) - \mathbb{E}J(\beta)| - t\mathbb{E}J(\beta) \right] \geq 0 \right\} \cap \mathcal{B}_{2s} \right) \\ \leq 2|N_j| \exp \left[-cn \cdot \min(t^2, t) \right]. \end{aligned}$$

907 Noting that $|N_j| \leq \left(\frac{ep}{s}\right)^s \left(\frac{3}{\varepsilon}\right)^s$ and $\gamma \leq 3\gamma_2$ completes the proof. \square

908 C.9 Proof of Lemma B.8

909 We will need the following lemma, whose proof is a straightforward calculation:

910 **Lemma C.6.** *Let D be a diagonal matrix and $A = (a_{ij}) = (I - L)(I - L)^T$ where L is a strictly*
 911 *lower triangular matrix. Then $\text{tr}(AD) \geq \text{tr}(D)$ with equality if and only if $A = I$ (i.e. $L = 0$).*

912 We now prove Lemma B.8. Write P for the permutation matrix corresponding to π . Now suppose
 913 that $B \in \mathbb{D}_p[\pi]$, so that $PBP^T = L$ and $P\tilde{B}(\pi)P^T = \tilde{L}$ are strictly lower triangular matrices. Then

$$\begin{aligned} \mathbb{E}[\ell(B)] &= \frac{1}{2} \text{tr} \left[(I - \tilde{B}(\pi))^{-1} (I - B) (I - B)^T (I - \tilde{B}(\pi))^{-T} \tilde{\Omega}(\pi) \right] \\ &= \frac{1}{2} \text{tr} \left[P(I - \tilde{B}(\pi))^{-1} (I - B) (I - B)^T (I - \tilde{B}(\pi))^{-T} \tilde{\Omega}(\pi) P^T \right] \\ &= \frac{1}{2} \text{tr} \left[(I - \tilde{L})^{-1} (I - L) (I - L)^T (I - \tilde{L})^{-T} P\tilde{\Omega}(\pi)P^T \right]. \end{aligned}$$

914 Note that $(I - \tilde{L})^{-1} (I - L)$ is lower triangular, so that $(I - \tilde{L})^{-1} (I - L) (I - L)^T (I - \tilde{L})^{-T} := A$ is
 915 of the form $A = (I - \bar{L})(I - \bar{L})^T$ for some strictly lower triangular matrix \bar{L} . In particular, restricted
 916 to $\mathbb{D}_p[\pi]$, $\mathbb{E}[\ell(B)]$ is of the form $\text{tr}(AD)$ for the diagonal matrix $D := P\tilde{\Omega}(\pi)P^T$. Inequality (36)
 917 then follows from Lemma C.6.

918 Finally, if there is equality in (36), then Lemma C.6 implies that $(I - \tilde{L})^{-1} (I - L) (I - L)^T (I - \tilde{L})^{-T}$
 919 $\tilde{L})^{-T} = I$, or

$$(I - L)(I - L)^T = (I - \tilde{L})^T (I - \tilde{L}),$$

920 and the desired claim follows from the uniqueness of the Cholesky decomposition.

921 C.10 Proof of Lemma B.9

922 Observe that for any $\pi \in \mathbb{S}_p$,

$$Q(\hat{B}) \leq Q(\tilde{B}(\pi)). \quad (58)$$

923 Moreover, we have the following alternative expression for Q :

$$Q(B) = \frac{1}{2n} \|\mathbf{X}(\tilde{B}(\hat{\pi}) - B) + \mathbf{E}(\hat{\pi})\|_F^2 + \rho_\lambda(B), \quad \text{for any } \hat{\pi} \in \hat{\Pi}. \quad (59)$$

924 Thus, using (58) and (59),

$$\begin{aligned} 0 &\leq Q(\tilde{B}(\pi)) - Q(\hat{B}) \\ &= \frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 - \frac{1}{2n} \|\mathbf{X}(\tilde{B}(\hat{\pi}) - \hat{B}) - \mathbf{E}(\hat{\pi})\|_F^2 + \rho_\lambda(\tilde{B}(\pi)) - \rho_\lambda(\hat{B}) \\ &= \frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 - \frac{1}{2n} \|\mathbf{E}(\hat{\pi})\|_F^2 - \frac{1}{2n} \|\mathbf{X}(\tilde{B}(\hat{\pi}) - \hat{B})\|_F^2 \\ &\quad + \frac{1}{n} \text{tr} \left(\mathbf{E}(\hat{\pi})^T \mathbf{X}(\tilde{B}(\hat{\pi}) - \hat{B}) \right) + \rho_\lambda(\tilde{B}(\pi)) - \rho_\lambda(\hat{B}). \end{aligned}$$

925 Since (58) holds for any π , this completes the proof.

926 D Auxiliary results

927 This section provides some additional results which are needed to prove Proposition B.10. This
 928 involves several steps: 1) Bounding the estimation error $\|\hat{B} - \tilde{B}(\hat{\pi})\|_r$ (Section D.1), 2) Controlling the
 929 terms (B.9b) and (B.9c) (Section D.2), and 3) Controlling (B.9a), which invokes the minimum-trace
 930 permutations Π_0 (Section D.3). After dealing with these prerequisites, we prove Proposition B.10 in
 931 Appendix D.4.

932 For any $\delta \in (0, 1)$, $\lambda \geq 0$, $\delta_0 > 0$, and $\pi \in \mathbb{S}_p$, define the following event:

$$\mathcal{G}(\delta_0, \lambda; \pi) = \left\{ \frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 - \frac{1}{2n} \|\mathbf{E}(\hat{\pi})\|_F^2 \leq \delta_0 \rho_\lambda(\tilde{B}(\pi)) \right\}. \quad (60)$$

933 The idea is to show that on this event (along with (74) and (75)), the desired conclusions hold. In
 934 Appendix D.3, we provide an explicit bound on the probability of $\mathcal{G}(\delta_0, \lambda; \pi)$.

935 D.1 Uniform deviation bounds

936 The purpose of this section is to control the estimation error $\|\hat{B} - \tilde{B}(\hat{\pi})\|_r$ via Proposition D.2, which
 937 is needed in the proof of Lemma D.7. This lemma—which is also proved in this Appendix—is a key
 938 prerequisite in the proof of Proposition B.10.

939 We start by establishing a general bound on the ℓ_r ($r = 1, 2$) estimation errors for a fixed design
 940 regression problem with a general regularizer ρ_λ . The objective here is to derive conditions under
 941 which we can guarantee such bounds for a fixed design problem, and then show that these conditions
 942 hold uniformly for all neighbourhood problems. The conditions we will need are familiar from the
 943 literature: A *Gaussian width condition* and a *restricted eigenvalue condition*.

944 For the rest of this subsection, we let $Z \in \mathbb{R}^{n \times m}$ and $w \in \mathbb{R}^n$ be a fixed matrix and fixed vector,
 945 respectively.

946 *Definition D.1* (Gaussian width). We say that the *Gaussian width (GW) condition* holds for (w, Z)
 947 relative to ρ_λ if there is a numerical constant $\delta \in (0, 1)$ such that

$$\frac{1}{n} |\langle w, Zu \rangle| \leq \delta \left[\frac{1}{2n} \|Zu\|_2^2 + \rho_\lambda(u) \right], \quad \forall u \in \mathbb{R}^m,$$

948 in which case we write $(w, Z) \in \text{GW}_{\rho_\lambda}(\delta)$. If this inequality is strict for all $u \neq 0$, we write
 949 $(w, Z) \in \text{GW}_{\rho_\lambda}^\circ(\delta)$.

950 We will be interested in the case where both w and Z are allowed to be random but independent.
 951 In this setting, for Gaussian designs considered in this paper, the GW condition holds with high
 952 probability for the ℓ_1 penalty (this follows from a standard Hölder inequality argument), and has

953 similarly been shown to hold for penalties induced by ℓ_q norms for $0 \leq q \leq 1$ [49]. Zhang and Zhang
 954 [74] provide a version of this condition that applies to general nonconvex regularizers.

955 Before we proceed, let us note the following key relation between model selection consistency and
 956 the GW condition:

957 **Lemma D.1.** *Consider the setup of Lemma B.3, namely, the regression problem $\mathbf{y} = Z\theta^* + \mathbf{w}$ but
 958 with $\theta^* = 0$. Then*

$$\mathcal{A}(\mathbf{w}/\delta, Z, 0)^c = \{(\mathbf{w}, Z) \in \text{GW}_{\rho_\lambda}^\circ(\delta)\}.$$

959 *Proof.* If $(\mathbf{w}, Z) \in \text{GW}_{\rho_\lambda}^\circ(\delta)$, then for any $u \neq 0$,

$$\begin{aligned} & \frac{\delta}{2n} \|Zu\|_2^2 - \frac{1}{n} \mathbf{w}^T Zu + \delta \rho_\lambda(u) > 0 \\ \iff & \frac{1}{2n} \|\mathbf{w}/\delta - Zu\|_2^2 + \rho_\lambda(u) > \frac{1}{2n} \|\mathbf{w}/\delta\|_2^2. \end{aligned}$$

960 The latter inequality implies

$$\{0\} = \arg \min_u \|\mathbf{w}/\delta - Zu\|_2^2 / (2n) + \rho_\lambda(u),$$

961 that is, 0 is the unique global minimizer of the right hand side. Recalling the definition of
 962 $\mathcal{A}(\mathbf{w}/\delta, Z, 0)$ in (22), we obtain the desired result. \square

963 Thus, in order to ensure the GW condition for (\mathbf{w}, Z_S) , it suffices to show that the corresponding
 964 regression problem is model selection consistent when the *true coefficients are all set to zero* and the
 965 noise variance is inflated by a factor of $1/\delta^2$. [74] refer to this property as *null-consistency*.

966 For any set $A \subset [m]$ and $\xi > 0$, define the following ‘‘cone’’:

$$C_{\rho_\lambda}(A, \xi) := \{u \in \mathbb{R}^m : \rho_\lambda(u_{A^c}) \leq \xi \rho_\lambda(u_A)\}. \quad (61)$$

967 This definition also depends on the ambient dimension m ; when we wish to emphasize this we will
 968 write $C_{\rho_\lambda}^m(A, \xi)$. The term ‘‘cone’’ here is used in an extended sense, in analogy with the ℓ_1 cone
 969 found in previous work.

970 **Definition D.2** (Generalized restricted eigenvalue). The *generalized restricted eigenvalue (RE)*
 971 *constant* of Z with respect to ρ_λ over a subset A is

$$\phi_{\rho_\lambda}^2(Z, A; \xi) := \inf \left\{ \frac{\|Zu\|_2^2}{n\|u\|_2^2} : u \in C_{\rho_\lambda}(A, \xi), u \neq 0 \right\}. \quad (62)$$

972 In the sequel, we often suppress the dependence of the generalized RE constants on λ and ξ , writing
 973 $\phi_\rho^2(Z, A) = \phi_{\rho_\lambda}^2(Z, A; \xi)$. Note that the usual restricted eigenvalue is equivalent to the special case
 974 $\rho_\lambda = \lambda \|\cdot\|_1$ [3].

975 Consider the usual linear regression set up, $y = Z\theta^* + w$, where $\theta^* \in \mathbb{R}^m$ and we define $S^* =$
 976 $\text{supp}(\theta^*)$. The following general result establishes that the two conditions $(w, Z) \in \text{GW}_\rho(\delta)$ and
 977 $\phi_\rho^2(Z, S^*) > 0$ are sufficient to bound the deviation $\hat{\theta} - \theta^*$:

978 **Theorem D.1.** *Assume $(w, Z) \in \text{GW}_{\rho_\lambda}(\delta)$ for some ρ_λ satisfying Condition A.1 and $\delta \in (0, 1)$. Let
 979 $\xi = \xi(\delta) := (1 + \delta)/(1 - \delta)$ and assume $\phi^2 := \phi_\rho^2(Z, S^*; \xi) > 0$. Then any $\hat{\theta} \in \hat{\Theta}_\lambda(Z\theta^* + w, Z)$
 980 satisfies*

$$\|\hat{\theta} - \theta^*\|_2 \leq C_2(\rho_\lambda, \xi, \phi) \cdot \|\theta^*\|_0^{1/2}, \quad (63)$$

$$\|\hat{\theta} - \theta^*\|_1 \leq C_1(\rho_\lambda, \xi, \phi) \cdot \|\theta^*\|_0. \quad (64)$$

981 **Remark D.1.** The constants in the previous theorem are given by

$$C_2(\rho_\lambda, \xi, \phi) = \frac{2\xi}{\phi^2} \lambda, \quad C_1(\rho_\lambda, \xi, \phi) = \frac{2\xi(1 + \xi)}{\phi^2} \lambda.$$

982 *Proof.* Recall that $S^* := \text{supp}(\theta^*)$. To lighten notation, for any vector u let $u_1 := u_{S^*}$, $u_2 := u_{(S^*)^c}$,
 983 and also $\Delta := \widehat{\theta} - \theta^*$. Then invoking the subadditivity of ρ_λ (this is a consequence of Condition A.1),

$$\begin{aligned} \rho_\lambda(\widehat{\theta}) - \rho_\lambda(\theta^*) &= \rho_\lambda(\Delta + \theta^*) - \rho_\lambda(\theta^*) \\ &= \rho_\lambda(\Delta_1 + \theta_1^*) + \rho_\lambda(\Delta_2) - \rho_\lambda(\theta_1^*) \\ &\geq -\rho_\lambda(\Delta_1) + \rho_\lambda(\Delta_2). \end{aligned} \quad (65)$$

984 It is straightforward to derive

$$\frac{1}{2n} \|y - Z\widehat{\theta}\|_2^2 - \frac{1}{2n} \|y - Z\theta^*\|_2^2 = \frac{1}{2n} \|Z\Delta\|_2^2 - \frac{1}{n} \langle w, Z\Delta \rangle. \quad (66)$$

985 Since $(w, Z) \in \text{GW}_{\rho_\lambda}(\delta)$, we can invoke the GW condition with $u = \Delta$,

$$-\frac{1}{n} \langle w, Z\Delta \rangle \geq -\frac{1}{n} |\langle w, Z\Delta \rangle| \geq -\delta \frac{1}{2n} \|Z\Delta\|_2^2 - \delta \rho_\lambda(\Delta). \quad (67)$$

986 It follows that

$$\begin{aligned} 0 &\geq \frac{1}{2n} \|y - Z\widehat{\theta}\|_2^2 - \frac{1}{2n} \|y - Z\theta^*\|_2^2 + \rho_\lambda(\widehat{\theta}) - \rho_\lambda(\theta^*) \\ &\geq \frac{1}{2n} \|Z\Delta\|_2^2 - \frac{1}{n} \langle w, Z\Delta \rangle - \rho_\lambda(\Delta_1) + \rho_\lambda(\Delta_2) \\ &\geq \frac{1-\delta}{2n} \|Z\Delta\|_2^2 - \delta \rho_\lambda(\Delta) - \rho_\lambda(\Delta_1) + \rho_\lambda(\Delta_2) \\ &= \frac{1-\delta}{2n} \|Z\Delta\|_2^2 - (1+\delta)\rho_\lambda(\Delta_1) + (1-\delta)\rho_\lambda(\Delta_2) \\ &= (1-\delta) \left[\frac{1}{2n} \|Z\Delta\|_2^2 + \rho_\lambda(\Delta_2) - \xi \rho_\lambda(\Delta_1) \right], \end{aligned} \quad (68)$$

987 where the first inequality by optimality of $\widehat{\theta}$, the second by (66), and the third by (67). The next
 988 line follows from an application of $\rho_\lambda(\Delta) = \rho_\lambda(\Delta_1) + \rho_\lambda(\Delta_2)$. Since $\delta < 1$ by assumption, it
 989 follows that $\rho_\lambda(\Delta_2) \leq \xi \rho_\lambda(\Delta_1)$ which implies $\Delta \in C_\rho(S^*, \xi(\delta))$.

990 Recalling the definition (62) of $\phi_\rho^2(Z, S^*)$, we conclude that $\frac{1}{2n} \|Z\Delta\|_2^2 \geq \frac{\phi^2}{2} \|\Delta\|_2^2$ which combined
 991 with (68), dropping $\rho_\lambda(\Delta_2)$, gives

$$0 \geq \frac{\phi^2}{2} \|\Delta\|_2^2 - \xi \rho_\lambda(\Delta_1).$$

992 Combining with the following (note $\|\Delta_1\|_0 \leq \|\theta^*\|_0$),

$$\rho_\lambda(\Delta_1) \leq \rho'_\lambda(0+) \|\Delta_1\|_1 \leq \rho'_\lambda(0+) \|\theta^*\|_0^{1/2} \|\Delta\|_2 \quad (69)$$

993 and re-arranging proves (63). For (64), since $\Delta \in C_\rho(S^*, \xi(\delta))$, we construct a set $M \subset [p]$ with
 994 $|M| = |S^*| = \|\theta^*\|_0$ such that $\Delta \in C_1(M, \xi(\delta))$. Then

$$\begin{aligned} \|\Delta\|_1 &= \|\Delta_M\|_1 + \|\Delta_{M^c}\|_1 \leq (1+\xi) \|\Delta_M\|_1 \\ &\leq (1+\xi) \|\theta^*\|_0^{1/2} \|\Delta_M\|_2 \\ &\leq \frac{2\xi(1+\xi)}{\phi^2} \cdot \rho'_\lambda(0+) \|\theta^*\|_0. \quad \square \end{aligned}$$

995 The GW condition is quantified by the constant $\delta \in (0, 1)$, and the restricted eigenvalue condition
 996 depends on the free parameter $\xi > 0$; these two are linked via the relation $\xi(\delta) = (1+\delta)/(1-\delta)$
 997 and play subtle roles in the proof. A slightly modified version of this result first appeared in Zhang
 998 and Zhang [74], under different assumptions. The particular version presented here is important to
 999 derive uniform bounds for all permutations, which we discuss next.

1000 In analogy with (23), define the following model selection exponent:

$$\psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta) := \inf_{0 \leq \sigma \leq \sigma_{\max}} \Phi_\lambda(\mathbf{X}, 0, \sigma^2/\delta^2). \quad (70)$$

1001 We often suppress the dependence on δ and write $\psi_\lambda(\mathbf{X}, \sigma_{\max}^2)$. Note that, in view of Lemma D.1,
 1002 $\psi_\lambda(\mathbf{X}, \sigma_{\max}^2)$ describes the conditional probability, given \mathbf{X} , that $(\sigma\mathbf{w}, \mathbf{X})$ violates a GW condition,
 1003 where $\mathbf{w} \sim \mathcal{N}_n(0, I_n)$ is independent of \mathbf{X} . More precisely,

$$\begin{aligned} \sup_{0 \leq \sigma \leq \sigma_{\max}} \mathbb{P}[(\sigma\mathbf{w}, \mathbf{X}) \notin \text{GW}_{\rho_\lambda}^\circ(\delta) \mid \mathbf{X}] &= \sup_{0 \leq \sigma \leq \sigma_{\max}} \exp[-\Phi_\lambda(\mathbf{X}, 0, \sigma^2/\delta^2)] \\ &= \exp[-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2)]. \end{aligned}$$

1004 We also recall the relation

$$\xi = \xi(\delta) = \frac{1 + \delta}{1 - \delta}. \quad (71)$$

1005 **Proposition D.2.** Assume that $\Sigma \succ 0$ and ρ_λ satisfies Condition A.1. Suppose $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$,
 1006 $\delta \in (0, 1)$, and define ξ by (71). Then there exist constants $c_0, c_1, c_2 > 0$ such that the following
 1007 holds: If

$$n > c_0 \frac{\sigma_{\max}^2 (1 + \xi)^2}{r_{\min}(\Sigma)} d \log p,$$

1008 then with probability at least $1 - c_1 \exp(-c_2 n) - p \binom{p}{d} \mathbb{E} \exp(-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta))$,

$$\|\widehat{\beta}_j(S) - \beta_j(S)\|_2 \leq C_2(\rho_\lambda, \xi, r_{\min}(\Sigma)) \cdot \|\beta_j(S)\|_0^{1/2}, \quad (72)$$

$$\|\widehat{\beta}_j(S) - \beta_j(S)\|_1 \leq C_1(\rho_\lambda, \xi, r_{\min}(\Sigma)) \cdot \|\beta_j(S)\|_0, \quad (73)$$

1009 uniformly over all $j \in [p]$ and $S \subset [p]_j$.

1010 For future reference, inspection of the proof shows that the conclusion of Proposition D.2 holds on
 1011 $\mathcal{E}(\delta, \lambda) \cap \mathcal{R}(\delta)$, where

$$\mathcal{E}(\delta, \lambda) = \left\{ (\tilde{\varepsilon}_j(S), \mathbf{X}_S) \in \text{GW}_{\rho_\lambda}^\circ(\delta), \forall j \in [p], S \subset [p]_j \right\}, \quad (74)$$

$$\mathcal{R}(\delta) = \left\{ \phi_\rho^2(\mathbf{X}_S, m_j(S)) \geq r_{\min}(\Sigma) > 0, \forall j \in [p], S \subset [p]_j \right\}. \quad (75)$$

1012 For regularizers that satisfy the lower bound in Condition A.1(c) we have the following control on
 1013 the exponent $\psi_\lambda(\mathbf{X}, \sigma_{\max}^2)$:

1014 **Proposition D.3.** Assume that $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$, and that ρ_λ satisfies Condition A.1(c). Then there
 1015 exist constants $c > 0$ and $C = C(\underline{\rho}_1, \underline{\rho}_0)$ such that for any $\delta \in (0, 1)$, if

$$\lambda \geq C \delta^{-1} \sigma_{\max} \|\Sigma\|^{1/4} \sqrt{\frac{(d+1) \log p}{n}} \quad (76)$$

1016 then $\mathbb{E} \exp(-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta)) \leq c \exp(-\min\{2(d+1) \log p, n\})$.

1017 The proof of Proposition D.3 follows from an argument similar to that in [74] and is omitted for
 1018 brevity. In order to prove Proposition D.2, we need the following two intermediate results, providing
 1019 uniform control on RE constants and GW conditions. Recall $\mathcal{E}(\delta, \lambda)$ as defined in (74).

1020 **Proposition D.4** (Uniform GW control). For any $\delta \in (0, 1)$ and $\lambda > 0$,

$$\mathbb{P}[\mathcal{E}(\delta, \lambda)] \geq 1 - p \binom{p}{d} \mathbb{E} \exp[-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta)].$$

1021 *Proof.* Fix $\delta \in (0, 1)$. By analogy with (31), for any neighbourhood $S \subset [p]_j$, let

$$\xi_j(S) := \Phi_\lambda(\mathbf{X}_S, 0, \omega_j^2(S)/\delta^2) \geq \psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta), \quad (77)$$

1022 where the inequality follows from (70) and $\omega_j^2(S) \leq \sigma_{\max}^2$. We follow the proof of Proposition B.5,
 1023 but with $\beta_j(S)$ replaced with 0, and $\tilde{\varepsilon}_j(S)$ replaced with $\tilde{\varepsilon}_j(S)/\delta$. To simplify, let $\mathcal{E} = \mathcal{E}(\delta, \lambda)$,

$$\mathcal{F}_S^j := \left\{ (\tilde{\varepsilon}_j(S), \mathbf{X}_S) \in \text{GW}_{\rho_\lambda}^\circ(\delta) \right\},$$

1024 and note that $\mathcal{E} = \bigcap_{j=1}^p \bigcap_{S \subset [p]_j} (\mathcal{F}_S^j)^c$. According to Lemma D.1, we have

$$\mathcal{F}_S^j = \mathcal{A}(\tilde{\varepsilon}_j(S)/\delta, \mathbf{X}, 0; S) = \mathcal{A}(\tilde{\varepsilon}_j(S)/\delta, \mathbf{X}_S, 0)$$

1025 where the second equality is by the same argument in the proof of Proposition B.5. Since $\tilde{\varepsilon}_j(S)/\delta \sim$
1026 $\mathcal{N}(0, [\omega_j^2(S)/\delta^2]I_n)$ independent of \mathbf{X}_S , we conclude, using Definition A.6, that

$$\mathbb{P}(\mathcal{F}_S^j \mid \mathbf{X}_S) = \exp[-\xi_j(S)],$$

1027 hence $\mathbb{P}(\mathcal{F}_S^j) \leq \mathbb{E} \exp[-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2)]$, $\forall S \subset [p]_j$, using the inequality in (77). The events \mathcal{F}_S^j are
1028 monotonic in S according to Corollary B.4. (The division of $\varepsilon_j(S)$ by δ does not change anything in
1029 that proof.) It follows that

$$\mathcal{E}^c = \bigcup_{j=1}^p \bigcup_{S \subset [p]_j} \mathcal{F}_S^j \subset \bigcup_{j=1}^p \bigcup_{T \in m_j(\Sigma)} \mathcal{F}_{M_j(T)}^j.$$

1030 Taking the union bound, and using $|m_j(\Sigma)| \leq \binom{p}{d}$ and

$$\mathbb{P}[\mathcal{F}_{M_j(T)}^j] \leq \mathbb{E} \exp[-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2)], \quad \forall T \in m_j(\Sigma),$$

1031 finishes the proof. \square

1032 **Proposition D.5** (Uniform RE control). *Assume $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$, $\Sigma \succ 0$, and ρ_λ satisfies Condi-*
1033 *tion A.1. There exist universal constants $c_0, c_1, c_2 > 0$, such that if*

$$n > c_0 \frac{\sigma_{\max}^2 (1 + \xi)^2}{r_{\min}(\Sigma)} d(\Sigma) \log p$$

1034 *then with probability at least $1 - c_1 \exp(-c_2 n)$,*

$$\inf_{1 \leq j \leq p} \inf_{S \subset [p]_j} \inf_{\substack{A \subset S \\ |A| \leq d}} \phi_\rho^2(\mathbf{X}_S, A; \xi) \geq r_{\min}(\Sigma).$$

1035 The proof of this proposition follows from the results in Raskutti et al. [48] and is omitted. Recalling
1036 the definition of $\mathcal{R}(\delta)$ in (75), combined with $m_j(S) = \|\beta_j(S)\|_0 \leq d$ (cf. Definition 4.1), Propo-
1037 sition D.5 implies that $\mathcal{R}(\delta)$ holds with probability at least $1 - c_1 \exp(-c_2 n)$. Let us show how
1038 Proposition D.2 follows.

1039 *Proof of Proposition D.2.* Recall the definitions of $\mathcal{E}(\delta, \lambda)$ in (74) and $\mathcal{R}(\delta)$ in (75). Propositions D.4
1040 and D.5 guarantee that

$$\mathbb{P}(\mathcal{R}(\delta) \cap \mathcal{E}(\delta, \lambda)) \geq 1 - c_1 \exp(-c_2 n) - p \binom{p}{d} \mathbb{E} \exp(-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta)).$$

1041 Thus, it suffices to deduce (72) and (73) whenever we are on the event $\mathcal{R}(\delta) \cap \mathcal{E}(\delta, \lambda)$. The case
1042 $\beta_j(S) = 0$ follows from Proposition D.4 and Lemma D.1, and the case $\beta_j(S) \neq 0$ follows from
1043 Theorem D.1 applied to the corresponding neighbourhood regression problems. \square

1044 D.2 Some intermediate lemmas

1045 Recall the definitions of $\mathcal{E}(\delta, \lambda)$ and $\mathcal{R}(\delta)$ in (74)–(75). We start with the following extension of GW
1046 bounds:

1047 **Lemma D.6.** *Let $\hat{\Delta} := \hat{B} - \tilde{B}(\hat{\pi})$. On $\mathcal{E}(\delta, \lambda)$, we have*

$$\frac{1}{n} \left| \text{tr} \left(\mathbf{E}(\hat{\pi})^T \mathbf{X} \hat{\Delta} \right) \right| < \delta \left[\frac{1}{2n} \|\mathbf{X} \hat{\Delta}\|_F^2 + \rho_\lambda(\hat{\Delta}) \right]. \quad (78)$$

1048 *Proof.* Let $\widehat{\Delta}_j := \widehat{\beta}_j - \widetilde{\beta}_j(\widehat{\pi})$ be the j th column of $\widehat{\Delta}$. Then

$$\frac{1}{n} \left| \text{tr} \left(\mathbf{E}(\widehat{\pi})^T \mathbf{X} \widehat{\Delta} \right) \right| \leq \frac{1}{n} \sum_{j=1}^p |\langle \widetilde{\epsilon}_j(\widehat{\pi}), \mathbf{X} \widehat{\Delta}_j \rangle|. \quad (79)$$

1049 According to (74), on $\mathcal{E}(\delta, \lambda)$, we have $(\widetilde{\epsilon}_j(S), \mathbf{X}_S) \in \text{GW}_{\rho_\lambda}^\circ(\delta)$ for all $S \subset [p]_j$. In particular,
1050 applying with $S = S_j(\widehat{\pi})$ and using $u = \widehat{\Delta}_j$ in the Definition D.1 of GW, we have

$$\frac{1}{n} |\langle \widetilde{\epsilon}_j(\widehat{\pi}), \mathbf{X} \widehat{\Delta}_j \rangle| < \delta \left[\frac{1}{2n} \|\mathbf{X} \widehat{\Delta}_j\|_2^2 + \rho_\lambda(\widehat{\Delta}_j) \right], \quad \forall j$$

1051 Summing over j and plugging into (79) yields (78). \square

1052 For any matrix $A = (a_{ij}) \in \mathbb{R}^{p \times p}$ and $S \subset [p] \times [p]$, let $A_{\langle S \rangle}$ denote the $p \times p$ matrix formed by
1053 zero-ing the elements outside of S , i.e.

$$(A_{\langle S \rangle})_{ij} = \begin{cases} a_{ij}, & (i, j) \in S, \\ 0, & (i, j) \notin S. \end{cases}$$

1054 In analogy with Condition 4.1 on signal strength, let us define

$$\tau_\lambda(\alpha; \Sigma) := \inf \left\{ \tau : \frac{\lambda^2}{\rho_\lambda(\tau)} \leq \frac{r_{\min}(\Sigma)}{\alpha} \right\} \quad (80)$$

1055 where we often suppress the dependence on Σ . Note that we can write Condition 4.1 equivalently as
1056 $\tau_* \geq \tau_\lambda(a_1)$.

1057 The next lemma is used to lower bound $\rho_\lambda(\widehat{B})$. The ℓ_0 case is easy to prove; for completeness we
1058 prove this for ℓ_1 and MCP.

1059 **Lemma D.7.** *Assume that ρ_λ satisfies Condition A.1, is right-differentiable with $\rho'_\lambda(0+) = \lambda$, and*

$$\tau_* \geq \tau_\lambda \left(\frac{2\xi}{1 - \delta_1} \right), \quad \text{for some } \delta_1 \in (0, 1) \quad (81)$$

1060 where $\xi = \xi(\delta)$ is defined by (71). Then, on $\mathcal{R}(\delta) \cap \mathcal{E}(\delta, \lambda)$,

$$\rho_\lambda(\widehat{B}) \geq \delta_1 \rho_\lambda(\widetilde{B}(\widehat{\pi})) + \rho_\lambda \left((\widehat{B} - \widetilde{B}(\widehat{\pi}))_{\langle \text{supp}(\widetilde{B}(\widehat{\pi}))^c \rangle} \right). \quad (82)$$

1061 *Proof.* To lighten the notation, let $\Delta = \widehat{B} - \widetilde{B}(\widehat{\pi})$, $S_1 = \text{supp}(\widetilde{B}(\widehat{\pi}))$, $\Delta_1 = \Delta_{\langle S_1 \rangle}$, and $\Delta_2 =$
1062 $\Delta_{\langle S_1^c \rangle}$. We have

$$\rho_\lambda(\Delta_1) \leq \lambda \|\Delta_1\|_1 \leq \lambda \|\widetilde{B}(\widehat{\pi})\|_0^{1/2} \|\Delta_1\|_2. \quad (83)$$

1063 Since we are on $\mathcal{R}(\delta) \cap \mathcal{E}(\delta, \lambda)$, Proposition D.2 yields the ℓ_2 deviation bound (72), which we use
1064 with $S = S_j(\widehat{\pi})$. Plugging into (83) and using $\|\Delta_1\|_2 \leq \|\Delta\|_2$,

$$\rho_\lambda(\Delta_1) \leq \lambda C_2(\rho_\lambda, \xi, r_{\min}(\Sigma)) \cdot \|\widetilde{B}(\widehat{\pi})\|_0. \quad (84)$$

1065 Trivially, we have $\rho_\lambda(\widetilde{B}(\widehat{\pi})) \geq \rho_\lambda(\tau_*) \|\widetilde{B}(\widehat{\pi})\|_0$, so that by (84)

$$\rho_\lambda(\Delta_1) \leq \frac{\lambda C_2(\rho_\lambda, \xi, r_{\min}(\Sigma))}{\rho_\lambda(\tau_*)} \cdot \rho_\lambda(\widetilde{B}(\widehat{\pi})) \leq (1 - \delta_1) \rho_\lambda(\widetilde{B}(\widehat{\pi})), \quad (85)$$

1066 where the last inequality follows from (81). Finally, note that

$$\begin{aligned} \rho_\lambda(\widehat{B}) &\geq \rho_\lambda(\widetilde{B}(\widehat{\pi})) + \rho_\lambda(\Delta_2) - \rho_\lambda(\Delta_1) \\ &\geq \delta_1 \rho_\lambda(\widetilde{B}(\widehat{\pi})) + \rho_\lambda(\Delta_2). \end{aligned}$$

1067 where the first inequality is by arguments similar to those leading to (65) and the second is by (85). \square

1068 *Remark D.2.* [62] use a slightly weaker beta-min condition in which only a constant fraction of the
 1069 edges of each DAG are assumed to be sufficiently large. Lemma D.7 and the ensuing arguments carry
 1070 through under such an assumption: Under Condition 3.5 in [62], we can use

$$\rho_\lambda(\tilde{B}(\hat{\pi})) \geq (1 - \eta_1)\rho_\lambda(\tau_*)\|\tilde{B}(\hat{\pi})\|_0,$$

1071 between (84) and (85) and obtain a bound similar to (82), with only the constants modified.

1072 The conclusion of Lemma D.7 is stronger than what we need in the sequel. We only use the weaker
 1073 inequality $\rho_\lambda(\tilde{B}) \geq \delta_1\rho_\lambda(\tilde{B}(\hat{\pi}))$ implied by (82).

1074 D.3 A bound on the sample residuals

1075 In this section, we prove the following result, which is used in the proof of Proposition B.10:

1076 **Proposition D.8.** *Assume $n > 4(C + 1)(d + 1)\log p$ for some $C > 0$ and let π_0 be a minimum-trace*
 1077 *permutation such that*

$$\frac{\rho_\lambda(\tilde{B}(\pi_0))}{\text{tr } \tilde{\Omega}(\pi_0)} \geq \frac{1}{\delta_0} \sqrt{\frac{50(C + 1)(d + 1)\log p}{n}}. \quad (86)$$

1078 Then for any $\delta_0 > 0$, $\mathbb{P}(\mathcal{G}(\delta_0, \lambda; \pi_0)) \geq 1 - 2e^{-C(d+1)\log p}$, i.e.

$$\mathbb{P}\left(\frac{1}{2n}\|\mathbf{E}(\pi_0)\|_F^2 - \frac{1}{2n}\|\mathbf{E}(\hat{\pi})\|_F^2 > \delta_0\rho_\lambda(\tilde{B}(\pi_0))\right) \leq 2e^{-C(d+1)\log p}.$$

1079 Define two functions by

$$h_n(u) := -\frac{u^2}{n} + \frac{2u}{\sqrt{n+1}} + \frac{1}{n+1}, \quad H_n(u) := \frac{u^2}{n} + \frac{2u}{\sqrt{n}}. \quad (87)$$

1080 These functions bound the deviations in the normed residuals $\tilde{\epsilon}_j(\pi)$, and will be used repeatedly in
 1081 the sequel. We note that

$$H_n(u) + h_n(u) \leq \frac{5u}{\sqrt{n}}, \quad u \geq n^{-1/2}. \quad (88)$$

1082 **Lemma D.9.** *Suppose $\mathbf{w} \sim \mathcal{N}_n(0, \sigma^2 I_n)$. Then for any $0 < u < n/\sqrt{n+1}$,*

$$\sigma^2(1 - h_n(u)) \leq \frac{1}{n}\|\mathbf{w}\|_2^2 \leq \sigma^2(1 + H_n(u)) \quad (89)$$

1083 with probability at least $1 - 2e^{-u^2/2}$.

1084 *Proof.* For $z \sim \mathcal{N}_n(0, I_n)$, we have the following useful bounds [see, e.g., 24, Corollary 1.2]:

$$\frac{n}{\sqrt{n+1}} \leq \mathbb{E}\|z\|_2 = \sqrt{2} \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})} \leq \sqrt{n}.$$

1085 Gaussian concentration implies that for any $u > 0$, both

$$\{\|\mathbf{w}\|_2 \leq \sigma(n/\sqrt{n+1} - u)\}, \quad \text{and} \quad \{\|\mathbf{w}\|_2 \geq \sigma(\sqrt{n} + u)\}$$

1086 hold with probability at most $e^{-u^2/2}$. Thus,

$$\mathbb{P}\left(\sigma^2\left(\frac{n}{\sqrt{n+1}} - u\right)^2 \leq \|\mathbf{w}\|_2^2 \leq \sigma^2(\sqrt{n} + u)^2\right) \geq 1 - 2e^{-u^2/2}. \quad (90)$$

1087 Re-writing (90) using (87) yields the desired result. \square

1088 **Lemma D.10.** Suppose $\mathbf{X} \stackrel{iid}{\sim} \mathcal{N}_p(0, \Sigma)$. Then for any $\pi \in \mathbb{S}_p$ and $0 < u < n/\sqrt{n+1}$,

$$\frac{1}{2} \operatorname{tr} \tilde{\Omega}(\pi) (1 - h_n(u)) \leq \frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 \leq \frac{1}{2} \operatorname{tr} \tilde{\Omega}(\pi) (1 + H_n(u)) \quad (91)$$

1089 with probability at least $1 - 2p \binom{p}{d} e^{-u^2/2}$.

1090 *Proof.* Note that for any $\pi \in \mathbb{S}_p$,

$$\frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 = \frac{1}{2n} \sum_{j=1}^p \|\tilde{\varepsilon}_j(\pi)\|_2^2 = \frac{1}{2n} \sum_{j=1}^p \|\tilde{\varepsilon}_j(S_j(\pi))\|_2^2. \quad (92)$$

1091 Thus it suffices to bound the deviations in $\|\tilde{\varepsilon}_j(S)\|_2$ for $S \subset [p]_j$. Consider the following events

$$\mathcal{G}_j(S) := \left\{ \frac{\omega_j^2(S)}{2} (1 - h_n(u)) \leq \frac{1}{2n} \|\tilde{\varepsilon}_j(S)\|_2^2 \leq \frac{\omega_j^2(S)}{2} (1 + H_n(u)) \right\}$$

1092 and let $\mathcal{G} := \bigcap_{j=1}^p \bigcap_{S \subset [p]_j} \mathcal{G}_j(S)$. By Lemma D.9, we have $\mathbb{P}(\mathcal{G}_j(S)) \geq 1 - 2e^{-u^2/2}$, for all

1093 $S \in [p]_j$. By a monotonicity argument (cf. (27)), we have $\mathcal{G} = \bigcap_{j=1}^p \bigcap_{S \in m_j(\Sigma)} \mathcal{G}_j(M_j(S))$.

1094 Applying the union bound and using (13),

$$\mathbb{P}(\mathcal{G}^c) \leq 2p \binom{p}{d} e^{-u^2/2}. \quad (93)$$

1095 Summing the inequalities defining $\mathcal{G}_j(S_j(\pi))$, over j , we conclude that (91) holds on \mathcal{G} . The proof is
1096 complete. \square

1097 Consider the (random) collection of permutations

$$\mathbb{S}_p^0 = \mathbb{S}_p^0(\delta_0; u) := \left\{ \pi \in \mathbb{S}_p : \frac{1}{2} \operatorname{tr} \tilde{\Omega}(\pi) [1 + H_n(u)] - \frac{1}{2} \operatorname{tr} \tilde{\Omega}(\hat{\pi}) [1 - h_n(u)] \leq \delta_0 \rho_\lambda(\tilde{B}(\pi)) \right\}.$$

1098 **Lemma D.11.** For any $\pi \in \mathbb{S}_p^0(\delta_0; u)$ and $0 < u < n/\sqrt{n+1}$, we have

$$\mathbb{P} \left(\frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 - \frac{1}{2n} \|\mathbf{E}(\hat{\pi})\|_F^2 > \delta_0 \rho_\lambda(\tilde{B}(\pi)) \right) \leq 2p \binom{p}{d} e^{-u^2/2}.$$

1099 *Proof.* Lemma D.10 implies that

$$\frac{1}{2n} \|\mathbf{E}(\pi)\|_F^2 - \frac{1}{2n} \|\mathbf{E}(\hat{\pi})\|_F^2 \leq \frac{1}{2} \operatorname{tr} \tilde{\Omega}(\pi) [1 + H_n(u)] - \frac{1}{2} \operatorname{tr} \tilde{\Omega}(\hat{\pi}) [1 - h_n(u)]$$

1100 with probability at least $1 - 2p \binom{p}{d} e^{-u^2/2}$. Since $\pi \in \mathbb{S}_p^0$, the right-side is bounded above by $\delta_0 \rho_\lambda \tilde{B}(\pi)$
1101 by definition, which establishes the claim. \square

1102 **Lemma D.12.** $1 - h_n(u) > 0$ for all $u \neq 0$, $n > 0$.

1103 *Proof.* Since $(u + \sqrt{n})^2 + 1 > 0$, re-writing this inequality yields

$$\begin{aligned} \frac{u^2}{n} + 1 &> \frac{2u}{\sqrt{n}} + \frac{1}{n} > \frac{2u}{\sqrt{n+1}} + \frac{1}{n+1} \\ \implies 1 + \frac{u^2}{n} - \frac{2u}{\sqrt{n+1}} - \frac{1}{n+1} &> 0 \end{aligned}$$

1104 Comparing with (87) yields the claim. \square

1105 *Proof of Proposition D.8.* Lemma D.11 implies that for a choice of
 1106 $u = \sqrt{2(C+1)(d+1)\log p}$, we have

$$\begin{aligned} \mathbb{P}\left(\frac{1}{2n}\|\mathbf{E}(\pi)\|_F^2 - \frac{1}{2n}\|\mathbf{E}(\hat{\pi})\|_F^2 > \delta_0\rho_\lambda(\tilde{B}(\pi))\right) &\leq 2p\binom{p}{d}e^{-(C+1)(d+1)\log p} \\ &\leq 2e^{-C(d+1)\log p} \end{aligned}$$

1107 for any $\pi \in \mathbb{S}_p^0$. Thus the claim will follow if we can show that $\pi_0 \in \mathbb{S}_p^0$. Note that

$$\begin{aligned} &\text{tr}\tilde{\Omega}(\pi_0)\left[1 + H_n(u)\right] - \text{tr}\tilde{\Omega}(\hat{\pi})\left[1 - h_n(u)\right] \\ &\stackrel{(i)}{\leq} \text{tr}\tilde{\Omega}(\pi_0)\left[H_n(u) + h_n(u)\right] \\ &\stackrel{(ii)}{\leq} \text{tr}\tilde{\Omega}(\pi_0)\sqrt{\frac{50(C+1)(d+1)\log p}{n}} \\ &\stackrel{(iii)}{\leq} \delta_0\rho_\lambda(\tilde{B}(\pi_0)), \end{aligned}$$

1108 where (i) follows from $\text{tr}\tilde{\Omega}(\pi_0) \leq \text{tr}\tilde{\Omega}(\hat{\pi})$ and Lemma D.12, (ii) follows by using (88) with
 1109 $u = \sqrt{2(C+1)(d+1)\log p}$, and (iii) follows from assumption (86). Hence, $\pi_0 \in \mathbb{S}_p^0$ and the proof
 1110 is complete. \square

1111 D.4 Proof of Proposition B.10

1112 *Proof.* Recall the definition of $\mathcal{G}(\delta_0, \lambda; \pi)$ in (60). Fix some π_0 such that $\tilde{B}(\pi_0) := \tilde{B}_{\min}$ satisfies
 1113 Condition 4.1 with $a_2 > 0$. Taking (arbitrarily) $C = 1$ and $\delta_0 = 10/a_2$ in Proposition D.8, we have

$$\mathbb{P}\left[\mathcal{G}(\delta_0, \lambda; \pi_0)^c\right] \leq 2e^{-(d+1)\log p}.$$

1114 Combined with Propositions D.5 and D.4, we obtain

$$\begin{aligned} &\mathbb{P}\left(\mathcal{G}(\delta_0, \lambda; \pi_0) \cap \mathcal{E}(\delta, \lambda) \cap \mathcal{R}(\delta)\right) \\ &\geq 1 - c_1 \exp(-c_2 \min\{n, (d+1)\log p\}) - p\binom{p}{d}\mathbb{E}\exp(-\psi_\lambda(\mathbf{X}, \sigma_{\max}^2; \delta)). \end{aligned}$$

1115 Thus, we may assume we are on $\mathcal{G}(\delta_0, \lambda; \pi_0) \cap \mathcal{E}(\delta, \lambda) \cap \mathcal{R}(\delta)$. Since we are on $\mathcal{E}(\delta, \lambda)$, we can
 1116 combine Lemma D.6 with Lemma B.9 (applied with $\pi = \pi_0$) to deduce (recall $\hat{\Delta} := \hat{B} - \tilde{B}(\hat{\pi})$)

$$\begin{aligned} \frac{1}{2n}\|\mathbf{X}\hat{\Delta}\|_F^2 + \rho_\lambda(\hat{B}) &\leq \frac{\delta}{2n}\|\mathbf{X}\hat{\Delta}\|_F^2 + \delta\rho_\lambda(\hat{\Delta}) \\ &\quad + \frac{1}{2n}\|\mathbf{E}(\pi_0)\|_F^2 - \frac{1}{2n}\|\mathbf{E}(\hat{\pi})\|_F^2 + \rho_\lambda(\tilde{B}(\pi_0)). \end{aligned}$$

1117 Dropping the prediction loss terms (those involving $\|\mathbf{X}\hat{\Delta}\|_F^2$), and using that we are on $\mathcal{G}(\delta_0, \lambda; \pi_0)$
 1118 to bound $\frac{1}{2n}\|\mathbf{E}(\pi_0)\|_F^2 - \frac{1}{2n}\|\mathbf{E}(\hat{\pi})\|_F^2$, we have after rearranging,

$$\begin{aligned} \rho_\lambda(\hat{B}) &\leq (1 + \delta_0)\rho_\lambda(\tilde{B}(\pi_0)) + \delta\rho_\lambda(\hat{B} - \tilde{B}(\hat{\pi})) \\ &\leq (1 + \delta_0)\rho_\lambda(\tilde{B}(\pi_0)) + \delta\rho_\lambda(\tilde{B}(\hat{\pi})) + \delta\rho_\lambda(\hat{B}). \end{aligned} \tag{94}$$

1119 Let $\delta_1 = 2\delta/(1 - \delta)$, so that $\xi/(1 - \delta_1) = (1 + \delta)/(1 - 3\delta)$ (cf. (71)). Furthermore, since $\delta < 1/3$
 1120 by assumption, $\delta_1 < 1$, so that Lemma D.7 implies $\rho_\lambda(\hat{B}) \geq \delta_1\rho_\lambda(\tilde{B}(\hat{\pi}))$ which gives (i) in (40).

1121 Since $\rho_\lambda(\tilde{B}(\hat{\pi})) \leq (1/\delta_1)\rho_\lambda(\hat{B})$, the bounds in (94) imply that

$$\rho_\lambda(\hat{B}) \leq (1 + \delta_0)\rho_\lambda(\tilde{B}(\pi_0)) + \frac{\delta}{\delta_1}\rho_\lambda(\hat{B}) + \delta\rho_\lambda(\hat{B}).$$

1122 Rearranging we get

$$\rho_\lambda(\hat{B}) \leq [1 - \delta(1 + \delta_1)/\delta_1]^{-1}(1 + \delta_0)\rho_\lambda(\tilde{B}(\pi_0)).$$

1123 We have $[1 - \delta(1 + \delta_1)/\delta_1]^{-1}(1 + \delta_0) = \frac{2}{1-\delta}(1 + \frac{10}{a_2})$, using $\delta_0 = 10/a_2$ and $\delta_1 = 2\delta/(1 - \delta)$ as
 1124 before. This proves (ii) in (40). \square