
Appendix: A Fast and Scalable Joint Estimator for Integrating Additional Knowledge in Learning Multiple Related Sparse Gaussian Graphical Models

Beilun Wang¹ Arshdeep Sekhon¹ Yanjun Qi¹

S:1 More about Method

Notations: $X_{n_i \times p}^{(i)}$ is the data matrix for the i -th task, which includes n_i data samples being described by p different feature variables. Then $n_{tot} = \sum_{i=1}^K n_i$ is the total number of data samples. We use notation $\Omega^{(i)}$ for the precision matrices and $\widehat{\Sigma}^{(i)}$ for the estimated covariance matrices. Given a p -dimensional vector $\mathbf{x} = (x_1, x_2, \dots, x_p)^T \in \mathbb{R}^p$, we denote the l_1 -norm of \mathbf{x} as $\|\mathbf{x}\|_1 = \sum_i |x_i|$. $\|\mathbf{x}\|_\infty = \max_i |x_i|$ is the l_∞ -norm of \mathbf{x} . Similarly, for a matrix X , let $\|X\|_1 = \sum_{i,j} |X_{i,j}|$ be the l_1 -norm of X and $\|X\|_\infty = \max_{i,j} |X_{i,j}|$ be the l_∞ -norm of X . $\|X\|_F = \sqrt{\sum_i \sum_j X_{i,j}^2}$

S:1.1 More about Solving JEEK

In Eq. (3.8), let $a_i = a_i^+ - a_i^-$ and $b = b^+ - b^-$. If $a_i \geq 0$, then $a_i^+ = a_i$ and $a_i^- = 0$. If $a_i < 0$, then $a_i^+ = 0$ and $a_i^- = -a_i$. The b^+ and b^- have the similar definition. Then Eq. (3.8) can be solved by the following small linear programming problem.

$$\begin{aligned} & \underset{a_i, b}{\operatorname{argmin}} \sum_i (w_i a_i^+ + w_i a_i^-) + K w_s b^+ + K w_s b^- \\ \text{Subject to: } & a_i^+ - a_i^- + b^+ - b^- \leq c_i + \frac{\lambda_n}{\min(w_i, w_s)}, \\ & a_i^+ - a_i^- + b^+ - b^- \geq c_i - \frac{\lambda_n}{\min(w_i, w_s)}, \\ & a_i^+, a_i^-, b^+, b^- \geq 0 \\ & i = 1, \dots, K \end{aligned}$$

¹Department of Computer Science, University of Virginia, Charlottesville, VA, USA. Correspondence to: Beilun Wang <bw4mw@virginia.edu>, Yanjun Qi <yanjun@virginia.edu>.

S:1.2 JEEK is Group entry-wise and parallelizing optimizable

JEEK can be easily paralleled. Essentially we just need to revise the “For loop” of step 6 and step 7 in Algorithm 1 into, for instance, “entry per machine” “entry per core”. Now We prove that JEEK is group entry-wise and parallelizing optimizable. We prove that our estimator can be optimized asynchronously in a group entry-wise manner.

Theorem S:1.1. (JEEK is Group entry-wise optimizable) Suppose we use JEEK to infer multiple inverse of covariance matrices summarized as $\widehat{\Omega}_{tot}$. $\{[\widehat{\Omega}_I^{(i)}]_{j,k}, [\widehat{\Omega}_S]_{j,k} | i = 1, \dots, K\}$ describes a group of $K + 1$ entries at (j, k) position. Varying $j \in \{1, 2, \dots, p\}$ and $k \in \{1, 2, \dots, p\}$, we have a total of $p \times p$ groups. If these groups are independently estimated by JEEK, then we have,

$$\bigcup_{j=1}^p \bigcup_{k=1}^p \{([\widehat{\Omega}_I^{(i)}]_{j,k} + [\widehat{\Omega}_S]_{j,k}) | i = 1, \dots, K\} = \widehat{\Omega}_{tot}. \quad (\text{S:1-1})$$

Proof. Eq. (3.8) are the small sub-linear programming problems on each group of entries. \square

S:2 Connecting to the Bayesian statistics

Our approach has a close connection to a hierarchical Bayesian model perspective. We show that the additional knowledge weight matrices are also the parameters of the prior distribution of $\Omega_I^{(i)}, \Omega_S$. In our formulation Eq. (3.7), $W_I^{(i)}, W_S$ are the additional knowledge weight matrices. From a hierarchical Bayesian view, the first level of the prior is a Gaussian distribution and the second level is a Laplace distribution. In the following section, we show that $W_I^{(i)}, W_S$ are also the parameters of Laplace distributions, which is a prior distribution of $\Omega_I^{(i)}, \Omega_S$.

Since by the definition, $\Omega_I^{(i)} \Omega_S = 0$. There are only two possible situations:

Case I ($\Omega_I^{(i)} \Omega_S = 0$):

$$X^{(i)} | \mu^{(i)}, \Omega^{(i)} \sim N(\mu^{(i)}, (\Omega^{(i)})^{-1}) \quad (\text{S:2-1})$$

$$\Omega_{j,k}^{(i)} | \mu^{(i)}, W_{I_{j,k}}^{(i)}, W_{S_{j,k}} = \Omega_{S_{j,k}} | \mu^{(i)}, W_{S_{j,k}} \quad (\text{S:2-2})$$

$$\begin{aligned} p(\Omega_{S_{j,k}} | \mu^{(i)}, W_{S_{j,k}}) \\ \propto e^{-(W_{S_{j,k}} | \Omega_{S_{j,k}})} \end{aligned} \quad (\text{S:2-3})$$

Here $\Omega_{S_{j,k}} | \mu^{(i)}, W_{S_{j,k}}$ follows a Laplace distribution with mean 0. $1/W_{S_{j,k}} > 0$ is the diversity parameter. The larger $W_{S_{j,k}}$ is, the distribution of $\Omega_{S_{j,k}} | \mu^{(i)}, W_{S_{j,k}}$ more likely concentrate on the 0. Namely, there will be the higher density for $\Omega_{S_{j,k}} = 0 | \mu^{(i)}, W_{S_{j,k}}$.

Case II ($\Omega_{S_{j,k}} = 0$):

$$X^{(i)} | \mu^{(i)}, \Omega^{(i)} \sim N(\mu^{(i)}, (\Omega^{(i)})^{-1}) \quad (\text{S:2-4})$$

$$\Omega_{j,k}^{(i)} | \mu^{(i)}, W_{I_{j,k}}^{(i)}, W_{S_{j,k}} = \Omega_{I_{j,k}}^{(i)} | \mu^{(i)}, W_{I_{j,k}}^{(i)} \quad (\text{S:2-5})$$

$$\begin{aligned} p(\Omega_{I_{j,k}}^{(i)} | \mu^{(i)}, W_{I_{j,k}}^{(i)}) \\ \propto e^{-(W_{I_{j,k}}^{(i)} | \Omega_{I_{j,k}}^{(i)})} \end{aligned} \quad (\text{S:2-6})$$

Here $\Omega_{I_{j,k}}^{(i)} | \mu^{(i)}, W_{I_{j,k}}^{(i)}$ follows a Laplace distribution with mean 0. $1/W_{I_{j,k}}^{(i)} > 0$ is the diversity parameter. The larger $W_{I_{j,k}}^{(i)}$ is, the distribution of $\Omega_{I_{j,k}}^{(i)} | \mu^{(i)}, W_{I_{j,k}}^{(i)}$ more likely concentrate on the 0. Namely, there will be the higher density for $\Omega_{I_{j,k}}^{(i)} = 0 | \mu^{(i)}, W_{I_{j,k}}^{(i)}$.

Therefore, we can combine the above two cases into the following one equation.

$$\begin{aligned} p(\Omega_{j,k}^{(i)} | \mu^{(i)}, W_{I_{j,k}}^{(i)}, W_{S_{j,k}}) \\ \propto e^{-(W_{I_{j,k}}^{(i)} | \Omega_{I_{j,k}}^{(i)} + W_{S_{j,k}} | \Omega_{S_{j,k}})} \end{aligned} \quad (\text{S:2-7})$$

Our final hierarchical Bayesian formulation consists of the Eq. (S:2-1) and Eq. (S:2-7). This model is a generalization of the model considered in the seminal paper on the Bayesian lasso (Park & Casella, 2008). The parameters $W_{I_{j,k}}^{(i)}, W_{S_{j,k}}$ in our general model are hyper-parameters that specify the shape of the prior distribution of each edges in $\Omega^{(i)}$. The negative log-posterior distribution of $\Omega^{(i)}$ is now given by:

$$\begin{aligned} -\log(\mathbb{P}(\Omega^{(i)} | X^{(i)}, \mu^{(i)}, W_{I_{j,k}}^{(i)}, W_{S_{j,k}})) \\ \propto -\log(\det(\Omega^{(i)})^{-1}) + < \Omega^{(i)}, \hat{\Sigma}^{(i)} > \\ + \sum_{j,k} (W_{I_{j,k}}^{(i)} | \Omega_{I_{j,k}}^{(i)}| + W_{S_{j,k}} | \Omega_{S_{j,k}}|) \end{aligned} \quad (\text{S:2-8})$$

Eq. (S:2-8) follows a weighted variation of Eq. (2.1).

S:3 More about Theoretical Analysis

S:3.1 Theorems and Proofs of three properties of kw-norm

In this sub-section, we prove the three properties of kw-norm used in Section 3.2. We then provide the convergence rate of our estimator based on these three properties.

- (i) kw-norm is a norm function if and only if any entries in W_I^{tot} and W_S^{tot} do not equal to 0.
- (ii) If the condition in (i) holds, kw-norm is a decomposable norm.
- (iii) If the condition in (i) holds, the dual norm of kw-norm is $\mathcal{R}^*(u) = \max(\|W_I^{tot} \circ u\|_\infty, \|W_S^{tot} \circ u\|_\infty)$.

S:3.1.1 NORM:

First we prove the correctness of the argument that kw-norm is a norm function by the following theorem:

Theorem S:3.1. *Eq. (3.6) is a norm if and only if $\forall 1 \leq j, k \leq p, W_{I_{j,k}}^{(i)} \neq 0$, and $W_{S_{j,k}} \neq 0$.*

This theorem gives the sufficient and necessary conditions to make kw-norm (Eq. (3.6)) a norm function.

S:3.1.2 DECOMPOSABLE NORM:

Then we show that kw-norm is a decomposable norm within a certain subspace. Before providing the theorem, we give the structural assumption of the parameter.

(IS-Sparsity): The 'true' parameter for Ω^{tot*} (multiple GGM structures) can be decomposed into two clear structures— Ω_I^{tot*} and Ω_S^{tot*} . Ω_I^{tot*} is exactly sparse with k_i non-zero entries indexed by a support set S_I and Ω_S^{tot*} is exactly sparse with k_s non-zero entries indexed by a support set S_S . $S_I \cap S_S = \emptyset$. All other elements equal to 0 (in $(S_I \cup S_S)^c$).

Definition S:3.2. (*IS-subspace*)

$$\mathcal{M}(S_I \cup S_S) = \{\theta_j = 0 | \forall j \notin S_I \cup S_S\} \quad (\text{S:3-1})$$

Theorem S:3.3. *Eq. (3.6) is a decomposable norm with respect to \mathcal{M} and \mathcal{M}^\perp*

S:3.1.3 DUAL NORM OF KW-NORM:

To obtain the final formulation Eq. (3.7) and its statistical convergence rate, we need to derive the dual norm formulation of kw-norm.

Theorem S:3.4. *The dual norm of kw-norm (Eq. (3.6)) is*

$$\mathcal{R}^*(u) = \max(\|W_I^{tot} \circ u\|_\infty, \|W_S^{tot} \circ u\|_\infty) \quad (\text{S:3-2})$$

The details of the proof are as follows.

S:3.1.4 PROOF OF THEOREM (S:3.1)

Lemma S:3.5. For kw -norm, $W_{I,j,k}^{tot} \neq 0$ and $W_{S,j,k}^{tot} \neq 0$ equals to $W_{I,j,k}^{tot} > 0$ and $W_{S,j,k}^{tot} > 0$.

Proof. If $W_{I,j,k}^{tot} < 0$, then $|W_{I,j,k}^{tot} \Omega_{I,j,k}^{tot}| = |W_{I,j,k}^{tot}| |\Omega_{I,j,k}^{tot}| = |-W_{I,j,k}^{tot}| |\Omega_{I,j,k}^{tot}|$. Notice that $-W_{I,j,k}^{tot} > 0$. \square

Proof. To prove the kw -norm is a norm, by Lemma (S:4.2) the only thing we need to prove is that $f(x) = \|W \circ x\|_1$ is a norm function if $W_{i,j} > 0$. 1. $f(ax) = \|aW \circ x\|_1 = |a| \|W \circ x\|_1 = |a| f(x)$. 2. $f(x+y) = \|W \circ (x+y)\|_1 = \|W \circ x + W \circ y\|_1 \leq \|W \circ x\|_1 + \|W \circ y\|_1 = f(x) + f(y)$. 3. $f(x) \geq 0$. 4. If $f(x) = 0$, then $\sum |W_{i,j} x_{i,j}| = 0$. Since $W_{i,j} \neq 0$, $x_{i,j} = 0$. Therefore, $x = 0$. Based on the above, $f(x)$ is a norm function. Since summation of norm is still a norm function, kw -norm is a norm function. \square

Futhurmore, we have the following Lemma:

Lemma S:3.6. The dual norm of $f(x)$ is $\|W \circ x\|_\infty$.

Proof. $f^*(u) = \sup_x \frac{\langle u, x \rangle}{\|W \circ x\|_1} = \sup_x \{ \langle u, x \rangle \mid \|W \circ x\|_1 \leq 1 \} = \|W \circ u\|_\infty$. \square

S:3.1.5 PROOF OF THEOREM (S:3.3)

Proof. Assume $u \in \mathcal{M}$ and $v \in \bar{\mathcal{M}}^\perp$, $\mathcal{R}(u+v) = \|W_I^{tot} \circ (u_I + v_I)\|_1 + \|W_S^{tot} \circ (u_S + v_S)\|_1 = \|W_I^{tot} \circ u_I\|_1 + \|W_S^{tot} \circ u_S\|_1 + \|W_I^{tot} \circ v_I\|_1 + \|W_S^{tot} \circ v_S\|_1 = \mathcal{R}(u) + \mathcal{R}(v)$. Therefore, kw -norm is a decomposable norm with respect to the subspace pair $(\mathcal{M}, \bar{\mathcal{M}}^\perp)$. \square

S:3.1.6 PROOF OF THEOREM (S:3.4)

Proof. Suppose $\mathcal{R}(\theta) = \sum_{\alpha \in I} c_\alpha \mathcal{R}_\alpha(\theta_\alpha)$, where $\sum_{\alpha \in I} \theta_\alpha = \theta$. Then the dual norm $\mathcal{R}^*(\cdot)$ can be derived by the following equation.

$$\begin{aligned} \mathcal{R}^*(u) &= \sup_\theta \frac{\langle \theta, u \rangle}{\mathcal{R}(\theta)} \\ &= \sup_{\theta_\alpha} \frac{\sum_\alpha \langle u, \theta_\alpha \rangle}{\sum_\alpha c_\alpha \mathcal{R}_\alpha(\theta_\alpha)} \\ &= \sup_{\theta_\alpha} \frac{\sum_\alpha \langle u/c_\alpha, \theta_\alpha \rangle}{\sum_\alpha \mathcal{R}_\alpha(\theta_\alpha)} \\ &\leq \sup_{\theta_\alpha} \frac{\sum_\alpha \mathcal{R}_\alpha^*(u/c_\alpha) \mathcal{R}_\alpha(\theta_\alpha)}{\sum_\alpha \mathcal{R}_\alpha(\theta_\alpha)} \\ &\leq \max_{\alpha \in I} \mathcal{R}_\alpha^*(u)/c_\alpha. \end{aligned} \quad (\text{S:3-3})$$

Therefore by Lemma (S:3.6), the dual norm of kw -norm is $\mathcal{R}^*(u) = \max(\|W_I^{tot} \circ u\|_\infty, \|W_S^{tot} \circ u\|_\infty)$. \square

S:3.2 Appendix: Proofs of Theorems about All Error Bounds of JEEK

S:3.2.1 DERIVATION OF THEOREM (4.1)

JEEK formulation Eq. (3.7) and EE-sGGM Eq. (2.5) are special cases of the following generic formulation:

$$\begin{aligned} &\underset{\theta}{\operatorname{argmin}} \mathcal{R}(\theta) \\ &\text{subject to: } \mathcal{R}^*(\theta - \hat{\theta}_n) \leq \lambda_n \end{aligned} \quad (\text{S:3-4})$$

Where $\mathcal{R}^*(\cdot)$ is the dual norm of $\mathcal{R}(\cdot)$,

$$\mathcal{R}^*(v) := \sup_{u \neq 0} \frac{\langle u, v \rangle}{\mathcal{R}(u)} = \sup_{\mathcal{R}(u) \leq 1} \langle u, v \rangle. \quad (\text{S:3-5})$$

Connecting Eq. (3.7) and Eq. (S:3-4), $\mathcal{R}(\cdot)$ is the kw -norm. $\hat{\theta}_n$ represents a close approximation of θ^* .

Following the unified framework (Negahban et al., 2009), we first decompose the parameter space into a subspace pair $(\mathcal{M}, \bar{\mathcal{M}}^\perp)$, where $\bar{\mathcal{M}}$ is the closure of \mathcal{M} . Here $\bar{\mathcal{M}}^\perp := \{v \in \mathbb{R}^p \mid \langle u, v \rangle = 0, \forall u \in \bar{\mathcal{M}}\}$. \mathcal{M} is the **model subspace** that typically has a much lower dimension than the original high-dimensional space. $\bar{\mathcal{M}}^\perp$ is the **perturbation subspace** of parameters. For further proofs, we assume the regularization function in Eq. (S:3-4) is **decomposable** w.r.t the subspace pair $(\mathcal{M}, \bar{\mathcal{M}}^\perp)$.

(C1) $\mathcal{R}(u+v) = \mathcal{R}(u) + \mathcal{R}(v)$, $\forall u \in \mathcal{M}, \forall v \in \bar{\mathcal{M}}^\perp$.

(Negahban et al., 2009) showed that most regularization norms are decomposable corresponding to a certain subspace pair.

Definition S:3.7. Subspace Compatibility Constant

Subspace compatibility constant is defined as $\Psi(\mathcal{M}, |\cdot|) := \sup_{u \in \mathcal{M} \setminus \{0\}} \frac{\mathcal{R}(u)}{|u|}$ which captures the relative value between the error norm $|\cdot|$ and the regularization function $\mathcal{R}(\cdot)$.

For simplicity, we assume there exists a true parameter θ^* which has the exact structure w.r.t a certain subspace pair. Concretely:

(C2) \exists a subspace pair $(\mathcal{M}, \bar{\mathcal{M}}^\perp)$ such that the true parameter satisfies $\operatorname{proj}_{\bar{\mathcal{M}}^\perp}(\theta^*) = 0$

Then we have the following theorem.

Theorem S:3.8. Suppose the regularization function in Eq. (S:3-4) satisfies condition (C1), the true parameter of Eq. (S:3-4) satisfies condition (C2), and λ_n satisfies that $\lambda_n \geq \mathcal{R}^*(\hat{\theta}_n - \theta^*)$. Then, the optimal solution $\hat{\theta}$ of Eq. (S:3-4) satisfies:

$$\mathcal{R}^*(\hat{\theta} - \theta^*) \leq 2\lambda_n \quad (\text{S:3-6})$$

$$\|\hat{\theta} - \theta^*\|_2 \leq 4\lambda_n \Psi(\bar{\mathcal{M}}) \quad (\text{S:3-7})$$

$$\mathcal{R}(\hat{\theta} - \theta^*) \leq 8\lambda_n \Psi(\bar{\mathcal{M}})^2 \quad (\text{S:3-8})$$

For the proposed JEEK model, $\mathcal{R}(\Omega^{tot}) = \|W_I^{tot} \circ \Omega_I^{tot}\|_1 + \|W_S^{tot} \circ \Omega_S^{tot}\|_1$. Based on the results in (Negahban et al., 2009), $\Psi(\bar{\mathcal{M}}) = \sqrt{k_i + k_s}$, where k_i and k_s are the total number of nonzero entries in Ω_I^{tot} and Ω_S^{tot} . Using $\mathcal{R}(\Omega^{tot}) = \|W_I^{tot} \circ \Omega_I^{tot}\|_1 + \|W_S^{tot} \circ \Omega_S^{tot}\|_1$ in Theorem (S:3.8), we have the following theorem (the same as Theorem (4.1)),

Theorem S:3.9. Suppose that $\mathcal{R}(\Omega^{tot}) = \|W_I^{tot} \circ \Omega_I^{tot}\|_1 + \|W_S^{tot} \circ \Omega_S^{tot}\|_1$ and the true parameter Ω^{tot*} satisfy the conditions (C1)(C2) and $\lambda_n \geq \mathcal{R}^*(\hat{\Omega}^{tot} - \Omega^{tot*})$, then the optimal point $\hat{\Omega}^{tot}$ of Eq. (3.7) has the following error bounds:

$$\begin{aligned} \max(\|W_I^{tot} \circ (\hat{\Omega}^{tot} - \Omega^{tot*})\|_\infty, \|W_S^{tot} \circ (\hat{\Omega}^{tot} - \Omega^{tot*})\|_\infty) &\leq 2\lambda_n \\ \|\hat{\Omega}^{tot} - \Omega^{tot*}\|_F &\leq 4\sqrt{k_i + k_s}\lambda_n \\ \|W_I^{tot} \circ (\hat{\Omega}_I^{tot} - \Omega_I^{tot*})\|_1 + \|W_S^{tot} \circ (\hat{\Omega}_S^{tot} - \Omega_S^{tot*})\|_1 &\leq 8(k_i + k_s)\lambda_n \end{aligned} \quad (\text{S:3-9})$$

S:3.2.2 PROOF OF THEOREM (S:3.8)

Proof. Let $\delta := \hat{\theta} - \theta^*$ be the error vector that we are interested in.

$$\begin{aligned} \mathcal{R}^*(\hat{\theta} - \theta^*) &= \mathcal{R}^*(\hat{\theta} - \hat{\theta}_n + \hat{\theta}_n - \theta^*) \\ &\leq \mathcal{R}^*(\hat{\theta}_n - \hat{\theta}) + \mathcal{R}^*(\hat{\theta}_n - \theta^*) \leq 2\lambda_n \end{aligned} \quad (\text{S:3-10})$$

By the fact that $\theta_{\mathcal{M}^\perp}^* = 0$, and the decomposability of \mathcal{R} with respect to $(\bar{\mathcal{M}}, \bar{\mathcal{M}}^\perp)$

$$\begin{aligned} \mathcal{R}(\theta^*) &= \mathcal{R}(\theta^*) + \mathcal{R}[\Pi_{\bar{\mathcal{M}}^\perp}(\delta)] - \mathcal{R}[\Pi_{\bar{\mathcal{M}}^\perp}(\delta)] \\ &= \mathcal{R}[\theta^* + \Pi_{\bar{\mathcal{M}}^\perp}(\delta)] - \mathcal{R}[\Pi_{\bar{\mathcal{M}}^\perp}(\delta)] \\ &\leq \mathcal{R}[\theta^* + \Pi_{\bar{\mathcal{M}}^\perp}(\delta) + \Pi_{\bar{\mathcal{M}}}(\delta)] + \mathcal{R}[\Pi_{\bar{\mathcal{M}}}(\delta)] \\ &\quad - \mathcal{R}[\Pi_{\bar{\mathcal{M}}^\perp}(\delta)] \\ &= \mathcal{R}[\theta^* + \delta] + \mathcal{R}[\Pi_{\bar{\mathcal{M}}}(\delta)] - \mathcal{R}[\Pi_{\bar{\mathcal{M}}^\perp}(\delta)] \end{aligned} \quad (\text{S:3-11})$$

Here, the inequality holds by the triangle inequality of norm. Since Eq. (S:3-4) minimizes $\mathcal{R}(\hat{\theta})$, we have $\mathcal{R}(\theta^* + \Delta) = \mathcal{R}(\hat{\theta}) \leq \mathcal{R}(\theta^*)$. Combining this inequality with Eq. (S:3-11), we have:

$$\mathcal{R}[\Pi_{\bar{\mathcal{M}}^\perp}(\delta)] \leq \mathcal{R}[\Pi_{\bar{\mathcal{M}}}(\delta)] \quad (\text{S:3-12})$$

Moreover, by Hlder's inequality and the decomposability of $\mathcal{R}(\cdot)$, we have:

$$\begin{aligned} \|\Delta\|_2^2 &= \langle \delta, \delta \rangle \leq \mathcal{R}^*(\delta) \mathcal{R}(\delta) \leq 2\lambda_n \mathcal{R}(\delta) \\ &= 2\lambda_n [\mathcal{R}(\Pi_{\bar{\mathcal{M}}}(\delta)) + \mathcal{R}(\Pi_{\bar{\mathcal{M}}^\perp}(\delta))] \leq 4\lambda_n \mathcal{R}(\Pi_{\bar{\mathcal{M}}}(\delta)) \\ &\leq 4\lambda_n \Psi(\bar{\mathcal{M}}) \|\Pi_{\bar{\mathcal{M}}}(\delta)\|_2 \end{aligned} \quad (\text{S:3-13})$$

where $\Psi(\bar{\mathcal{M}})$ is a simple notation for $\Psi(\bar{\mathcal{M}}, \|\cdot\|_2)$.

Since the projection operator is defined in terms of $\|\cdot\|_2$ norm, it is non-expansive: $\|\Pi_{\bar{\mathcal{M}}}(\Delta)\|_2 \leq \|\Delta\|_2$. Therefore, by Eq. (S:3-13), we have:

$$\|\Pi_{\bar{\mathcal{M}}}(\delta)\|_2 \leq 4\lambda_n \Psi(\bar{\mathcal{M}}), \quad (\text{S:3-14})$$

and plugging it back to Eq. (S:3-13) yields the error bound Eq. (S:3-7).

Finally, Eq. (S:3-8) is straightforward from Eq. (S:3-12) and Eq. (S:3-14).

$$\begin{aligned} \mathcal{R}(\delta) &\leq 2\mathcal{R}(\Pi_{\bar{\mathcal{M}}}(\delta)) \\ &\leq 2\Psi(\bar{\mathcal{M}}) \|\Pi_{\bar{\mathcal{M}}}(\delta)\|_2 \leq 8\lambda_n \Psi(\bar{\mathcal{M}})^2. \end{aligned} \quad (\text{S:3-15})$$

□

S:3.2.3 CONDITIONS OF PROVING ERROR BOUNDS OF JEEK

JEEK achieves similar convergence rates as the SIMULE (Wang et al., 2017b) (W-SIMULE with no additional knowledge) and FASJEM estimator (Wang et al., 2017a). The other multiple sGGMs estimation methods have not provided such convergence rate analysis.

To derive the statistical error bound of JEEK, we need to assume that $\text{inv}(T_v(\hat{\Sigma}^{tot}))$ are well-defined. This is ensured by assuming that the true $\Omega^{(i)*}$ satisfy the following conditions (Yang et al., 2014):

(C-MinInf- Σ): The true $\Omega^{(i)*}$ Eq. (3.7) have bounded induced operator norm, i.e., $\|\Omega^{(i)*}\|_\infty := \sup_{w \neq 0 \in \mathbb{R}^p} \frac{\|\Sigma^{(i)*} w\|_\infty}{\|w\|_\infty} \leq \kappa_1$.

(C-Sparse- Σ): The true covariance matrices $\Sigma^{(i)*}$ are ‘‘approximately sparse’’ (following (Bickel & Levina, 2008)). For some constant $0 \leq q < 1$ and $c_0(p)$, $\max_i \sum_{j=1}^p |[\Sigma^{(i)*}]_{ij}|^q \leq c_0(p)$.¹

¹This indicates for some positive constant d , $[\Sigma^{(i)*}]_{jj} \leq d$ for all diagonal entries. Moreover, if $q = 0$, then this condition reduces to $\Sigma^{(i)*}$.

We additionally require $\inf_{w \neq 0 \in \mathbb{R}^p} \frac{\|\Omega^{(i)*} w\|_\infty}{\|w\|_\infty} \geq \kappa_2$.

S:3.2.4 PROOF OF COROLLARY (4.2)

Proof. In the following proof, we re-denote the following

two notations: $\Sigma_{tot} := \begin{pmatrix} \Sigma^{(1)} & 0 & \dots & 0 \\ 0 & \Sigma^{(2)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Sigma^{(K)} \end{pmatrix}$

and

$$\Omega_{tot} := \begin{pmatrix} \Omega^{(1)} & 0 & \dots & 0 \\ 0 & \Omega^{(2)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Omega^{(K)} \end{pmatrix}$$

The condition (C-Sparse Σ) and condition (C-MinInf Σ) also hold for Ω_{tot}^* and Σ_{tot}^* . In order to utilize Theorem (S:3.9) for this specific case, we only need to show that $\|\Omega_{tot}^* - [T_\nu(\hat{\Sigma}_{tot})]^{-1}\|_\infty \leq \lambda_n$ for the setting of λ_n in the statement:

$$\begin{aligned} & \|\Omega_{tot}^* - [T_\nu(\hat{\Sigma}_{tot})]^{-1}\|_\infty \\ &= \|[T_\nu(\hat{\Sigma}_{tot})]^{-1}(T_\nu(\hat{\Sigma}_{tot})\Omega_{tot}^* - I)\|_\infty \\ &\leq \| [T_\nu(\hat{\Sigma}_{tot})w] \|_\infty \|T_\nu(\hat{\Sigma}_{tot})\Omega_{tot}^* - I\|_\infty \\ &= \| [T_\nu(\hat{\Sigma}_{tot})]^{-1} \|_\infty \|\Omega_{tot}^*(T_\nu(\hat{\Sigma}_{tot}) - \Sigma_{tot}^*)\|_\infty \\ &\leq \| [T_\nu(\hat{\Sigma}_{tot})]^{-1} \|_\infty \|\Omega_{tot}^*\|_\infty \|T_\nu(\hat{\Sigma}_{tot}) - \Sigma_{tot}^*\|_\infty. \end{aligned} \quad (\text{S:3-16})$$

We first compute the upper bound of $\| [T_\nu(\hat{\Sigma}_{tot})]^{-1} \|_\infty$. By the selection ν in the statement, Lemma (S:4.2) and Lemma (S:4.3) hold with probability at least $1 - 4/p'^{\tau-2}$. Armed with Eq. (S:4-9), we use the triangle inequality of norm and the condition (C-Sparse Σ): for any w ,

$$\begin{aligned} \|T_\nu(\hat{\Sigma}_{tot})w\|_\infty &= \|T_\nu(\hat{\Sigma}_{tot})w - \Sigma w + \Sigma w\|_\infty \\ &\geq \|\Sigma w\|_\infty - \|(T_\nu(\hat{\Sigma}_{tot}) - \Sigma)w\|_\infty \\ &\geq \kappa_2 \|w\|_\infty - \|(T_\nu(\hat{\Sigma}_{tot}) - \Sigma)w\|_\infty \\ &\geq (\kappa_2 - \|(T_\nu(\hat{\Sigma}_{tot}) - \Sigma)w\|_\infty) \|w\|_\infty \end{aligned} \quad (\text{S:3-17})$$

Where the second inequality uses the condition (C-Sparse Σ). Now, by Lemma (S:4.2) with the selection of ν , we have

$$\|T_\nu(\hat{\Sigma}_{tot}) - \Sigma\|_\infty \leq c_1 \left(\frac{\log(Kp')}{n_{tot}} \right)^{(1-q)/2} c_0(p) \quad (\text{S:3-18})$$

where c_1 is a constant related only on τ and $\max_i \Sigma_{ii}$. Specifically, it is defined as $6.5(16(\max_i \Sigma_{ii})\sqrt{10\tau})^{1-q}$. Hence, as long as $n_{tot} > (\frac{2c_1 c_0(p)}{\kappa_2})^{\frac{2}{1-q}} \log p'$ as stated, so that $\|T_\nu(\hat{\Sigma}_{tot}) - \Sigma\|_\infty \leq \frac{\kappa_2}{2}$, we can conclude that $\|T_\nu(\hat{\Sigma}_{tot})w\|_\infty \geq \frac{\kappa_2}{2} \|w\|_\infty$, which implies $\| [T_\nu(\hat{\Sigma}_{tot})]^{-1} \|_\infty \leq \frac{2}{\kappa_2}$.

The remaining term in Eq. (S:3-16) is $\|T_\nu(\hat{\Sigma}_{tot}) - \Sigma_{tot}^*\|_\infty$; $\|T_\nu(\hat{\Sigma}_{tot}) - \Sigma_{tot}^*\|_\infty \leq \|T_\nu(\hat{\Sigma}_{tot}) - \hat{\Sigma}_{tot}\|_\infty + \|\hat{\Sigma}_{tot} - \Sigma_{tot}^*\|_\infty$. By construction of $T_\nu(\cdot)$ in (C-Threshold) and by Lemma (S:4.3), we can confirm that $\|T_\nu(\hat{\Sigma}_{tot}) - \hat{\Sigma}_{tot}\|_\infty$ as well as $\|\hat{\Sigma}_{tot} - \Sigma_{tot}^*\|_\infty$ can be upper-bounded by ν .

Therefore,

$$\begin{aligned} & \max(\|W_I^{tot} \circ (\Omega^{tot*} - \text{inv}(T_\nu(\hat{\Sigma}^{tot})))\|_\infty, \\ & \|W_S^{tot} \circ (\Omega^{tot*} - \text{inv}(T_\nu(\hat{\Sigma}^{tot})))\|_\infty) \\ & \leq O(\max_{j,k} \max(W_I^{tot}_{j,k}, W_S^{tot}_{j,k}) \sqrt{\frac{\log(Kp)}{n_{tot}}}) \end{aligned} \quad (\text{S:3-19})$$

By combining all together, we can confirm that the selection of λ_n satisfies the requirement of Theorem (S:3.9), which completes the proof. \square

S:4 Appendix: More Background of Proxy Backward mapping and Theorems of T_ν Being Invertible

The first row of Figure 1 summarizes the EE-sGGMs. Two important concepts:

(1) Backward Mapping: The Gaussian distribution is naturally an exponential-family distribution. Based on (Wainwright & Jordan, 2008), learning an exponential family distribution from data means to estimate its canonical parameter. For an exponential family distribution, computing the canonical parameter through vanilla graphical model MLE can be expressed as a backward mapping (the first step in Figure 1). For a Gaussian, the backward mapping is easily computable as the inverse of the sample covariance matrix. More details in Section (S:4.1).

(2) Proxy Backward Mapping: When being high-dimensional, we can not compute the backward mapping of Gaussian through the inverse of the sample covariance matrix. Now the key is to find a closed-form and statistically guaranteed estimator as the proxy backward mapping under high-dimensional cases. By the conclusion given by the EE-sGGM, we choose $\{([T_\nu(\hat{\Sigma}^{(i)})]^{-1})\}$ as the proxy backward mapping for $\{\Omega^{(i)}\}$.

$$[T_\nu(A)]_{ij} := \rho_v(A_{ij}) \quad (\text{S:4-1})$$

where $\rho_v(\cdot)$ is chosen to be a soft-thresholding function.

S:4.1 More About Background: backward mapping for an exponential-family distribution:

The solution of vanilla graphical model MLE can be expressed as a backward mapping (Wainwright & Jordan, 2008) for an exponential family distribution. It estimates the model parameters (canonical parameter θ) from certain (sample) moments. We provide detailed explanations about backward mapping of exponential families, backward mapping for Gaussian special case and backward mapping for differential network of GGM in this section.

Backward mapping: Essentially the vanilla graphical model MLE can be expressed as a backward mapping that computes the model parameters corresponding to some given moments in an exponential family distribution. For instance, in the case of learning GGM with vanilla MLE, the backward mapping is $\hat{\Sigma}^{-1}$ that estimates Ω from the sample covariance (moment) $\hat{\Sigma}$.

Suppose a random variable $X \in \mathbb{R}^p$ follows the exponential family distribution:

$$\mathbb{P}(X; \theta) = h(X) \exp\{\langle \theta, \phi(X) \rangle - A(\theta)\} \quad (\text{S:4-2})$$

Where $\theta \in \Theta \subset \mathbb{R}^d$ is the canonical parameter to be estimated and Θ denotes the parameter space. $\phi(X)$ denotes the sufficient statistics as a feature mapping function $\phi: \mathbb{R}^p \rightarrow \mathbb{R}^d$, and $A(\theta)$ is the log-partition function. We then define mean parameters v as the expectation of $\phi(X)$: $v(\theta) := \mathbb{E}[\phi(X)]$, which can be the first and second moments of the sufficient statistics $\phi(X)$ under the exponential family distribution. The set of all possible moments by the moment polytope:

$$\mathcal{M} = \{v | \exists p \text{ is a distribution s.t. } \mathbb{E}_p[\phi(X)] = v\} \quad (\text{S:4-3})$$

Mostly, the graphical model inference involves the task of computing moments $v(\theta) \in \mathcal{M}$ given the canonical parameters $\theta \in \textcircled{\text{H}}$. We denote this computing as **forward mapping**:

$$\mathcal{A}: \textcircled{\text{H}} \rightarrow \mathcal{M} \quad (\text{S:4-4})$$

The learning/estimation of graphical models involves the task of the reverse computing of the forward mapping, the so-called **backward mapping** (Wainwright & Jordan, 2008). We denote the interior of \mathcal{M} as \mathcal{M}^0 . **backward mapping** is defined as:

$$\mathcal{A}^*: \mathcal{M}^0 \rightarrow \textcircled{\text{H}} \quad (\text{S:4-5})$$

which does not need to be unique. For the exponential family distribution,

$$\mathcal{A}^*: v(\theta) \rightarrow \theta = \nabla A^*(v(\theta)). \quad (\text{S:4-6})$$

Where $\mathcal{A}^*(v(\theta)) = \sup_{\theta \in \textcircled{\text{H}}} \langle \theta, v(\theta) \rangle - A(\theta)$.

Backward Mapping: Gaussian Case If a random variable $X \in \mathbb{R}^p$ follows the Gaussian Distribution $N(\mu, \Sigma)$, then $\theta = (\Sigma^{-1}\mu, -\frac{1}{2}\Sigma^{-1})$. The sufficient statistics $\phi(X) = (X, XX^T)$, $h(x) = (2\pi)^{-\frac{k}{2}}$, and the log-partition function

$$A(\theta) = \frac{1}{2}\mu^T \Sigma^{-1} \mu + \frac{1}{2} \log(|\Sigma|) \quad (\text{S:4-7})$$

When performing the inference of Gaussian Graphical Models, it is easy to estimate the mean vector $v(\theta)$, since it equals to $\mathbb{E}[X, XX^T]$.

When learning the GGM, we estimate its canonical parameter θ through vanilla MLE. Because Σ^{-1} is one entry of θ we can use the backward mapping to estimate Σ^{-1} .

$$\begin{aligned} \theta &= (\Sigma^{-1}\mu, -\frac{1}{2}\Sigma^{-1}) = \mathcal{A}^*(v) = \nabla A^*(v) \\ &= ((\mathbb{E}_\theta[XX^T] - \mathbb{E}_\theta[X]\mathbb{E}_\theta[X]^T)^{-1}\mathbb{E}_\theta[X], \\ &\quad -\frac{1}{2}(\mathbb{E}_\theta[XX^T] - \mathbb{E}_\theta[X]\mathbb{E}_\theta[X]^T)^{-1}). \end{aligned} \quad (\text{S:4-8})$$

By plugging in Eq. (S:4-7) into Eq. (S:4-6), we get the backward mapping of Ω as $(\mathbb{E}_\theta[XX^T] - \mathbb{E}_\theta[X]\mathbb{E}_\theta[X]^T)^{-1} = \hat{\Sigma}^{-1}$, easily computable from the sample covariance matrix.

S:4.2 Theorems of T_v Being Invertible

Based on (Yang et al., 2014) for any matrix A , the element wise operator T_v is defined as:

$$[T_v(A)]_{ij} = \begin{cases} A_{ii} + v & \text{if } i = j \\ \text{sign}(A_{ij})(|A_{ij}| - v) & \text{otherwise, } i \neq j \end{cases}$$

Suppose we apply this operator T_v to the sample covariance matrix $\frac{X^T X}{n}$ to obtain $T_v(\frac{X^T X}{n})$. Then, $T_v(\frac{X^T X}{n})$ under high dimensional settings will be invertible with high probability, under the following conditions:

Condition-1 (Σ -Gaussian ensemble) Each row of the design matrix $X \in \mathbb{R}^{n \times p}$ is i.i.d sampled from $N(0, \Sigma)$.

Condition-2 The covariance Σ of the Σ -Gaussian ensemble is strictly diagonally dominant: for all row i , $\delta_i := \Sigma_{ii} - \sum_{j \neq i} \Sigma_{ij} \geq \delta_{\min} > 0$ where δ_{\min} is a large enough constant so that $\|\Sigma\|_\infty \leq \frac{1}{\delta_{\min}}$.

This assumption guarantees that the matrix $T_v(\frac{X^T X}{n})$ is invertible, and its induced ℓ_∞ norm is well bounded. Then the following theorem holds:

Theorem S:4.1. Suppose Condition-1 and Condition-2 hold. Then for any $v \geq 8(\max_i \Sigma_{ii})\sqrt{(\frac{10\tau \log p'}{n})}$, the matrix $T_v(\frac{X^T X}{n})$ is invertible with probability at least $1 - 4/p'^{\tau-2}$ for $p' := \max\{n, p\}$ and any constant $\tau > 2$.

Then we provide the error bound of T_v in the first lemma of Section (S:4.3) and use it in deriving the error bound of JEEK.

S:4.3 Useful lemma(s) of Error Bounds of (Proxy) Backward Mapping

Lemma S:4.2. (Theorem 1 of (Rothman et al., 2009)). Let δ be $\max_{ij} |\left[\frac{X^T X}{n}\right]_{ij} - \Sigma_{ij}|$. Suppose that $\nu > 2\delta$. Then, under the conditions (C-Sparse Σ), and as $\rho_v(\cdot)$ is a soft-threshold function, we can deterministically guarantee that the spectral norm of error is bounded as follows:

$$\|T_v(\hat{\Sigma}) - \Sigma\|_\infty \leq 5\nu^{1-q}c_0(p) + 3\nu^{-q}c_0(p)\delta \quad (\text{S:4-9})$$

Lemma S:4.3. (Lemma 1 of (Ravikumar et al., 2011)). Let \mathcal{A} be the event that

$$\left\|\frac{X^T X}{n} - \Sigma\right\|_\infty \leq 8(\max_i \Sigma_{ii})\sqrt{\frac{10\tau \log p'}{n}} \quad (\text{S:4-10})$$

where $p' := \max(n, p)$ and τ is any constant greater than 2. Suppose that the design matrix X is i.i.d. sampled from Σ -Gaussian ensemble with $n \geq 40 \max_i \Sigma_{ii}$. Then, the probability of event \mathcal{A} occurring is at least $1 - 4/p'^{\tau-2}$.

S:5 Design W_S and $W_I^{(i)}$: connections with related work and real-world applications

In this section, we showcase with specific examples that our proposed model JEEK can easily incorporate edge-level (like distance) as well as node-based (like hubs or groups) knowledge for the joint estimation of multiple graphs. To this end, we introduce four different choices of W_S^{tot} and W_I^{tot} in our formulation Eq. (3.7). By simply designing different choices of W_S^{tot} and W_I^{tot} , we can express different kinds of additional knowledge explicitly without changing the optimization algorithm.

Specifically, we design W_S and $W_I^{(i)}$ for cases like:

- (1) the additional knowledge is available in the form of a $p \times p$ matrix W . For instance distance matrix among brain regions in neuroscience study belongs to this type;
- (2) the existing knowledge is not in the form of matrix about nodes. We need to design W for such cases, for example the information of known hub nodes or the information of how nodes fall into groups (e.g., genes belonging to the same pathway or locations).

For the second kind, we showcase three different designs of weight matrices for representing (a) known co-Hub nodes, (b) perturbed hub nodes, and (c) node grouping information.

The design of knowledge matrices is loosely related to the different structural assumptions used by the JGL studies as ((Mohan et al., 2013), (Danaher et al., 2013)). For example, JGL can use specially designed norms like the one proposed in (Mohan et al., 2013) to push multiple graphs to have a similar set of nodes as hubs. However JGL can not model additional knowledge like a specific set of nodes are hub nodes (like we know node j is a hub node). Differently, JEEK can design $\{W_I^{(i)}, W_S\}$ for incorporating such knowledge. Essentially JEEK is complementary to JGL because they capture different type of prior information.

S:5.1 Case study I: Knowledge as matrix form like a distance matrix or some known edges

The first example we consider is exploiting a spatial prior to jointly estimate brain connectivity for multiple subject groups. Over time, neuroscientists have gathered considerable knowledge regarding the spatial and anatomical information underlying brain connectivity (i.e. short edges and certain anatomical regions are more likely to be connected (Watts & Strogatz, 1998)). Previous studies enforce these priors via a matrix of weights, W , corresponding to edges. To use our proposed model JEEK for such tasks, we can similarly choose $W = W_I^{(i)} = W_S$ in Eq. (3.7).

S:5.2 Case study II: Knowledge of co-hub nodes

The structure assumption we consider is graphs with co-hub nodes. Namely, there exists a set of nodes $NId = \{j | j \in \{1, 2, \dots, p\}\}$ such that $\Omega_{j,k}^{(i)} \neq 0, \forall i \in \{1, 2, \dots, K\}$ and $k \in \{1, \dots, p\}$. The above sub-figure of Figure S:4 is an example of the co-hub nodes.

A so-called JGL-hub (Mohan et al., 2013) estimator chooses $\mathcal{R}'(\cdot) = \sum_{i < i'} P_q(\Omega^{(i)} - \Omega^{(i')})$ in Eq. (5.2) to account for the co-hub structure assumption. Here $P_q(\Theta_1, \Theta_2, \dots, \Theta_k) = 1/2 \|\Theta_1, \dots, \Theta_k\|_{\ell_1, \ell_q}$. Θ_i is a symmetric matrix and $\|\cdot\|_{\ell_1, \ell_q}$ is the notation of ℓ_1, ℓ_q -norm. JGL-hub formulation needs a complicated ADMM solution with computationally expensive SVD steps.

We design W_S and $W_I^{(i)}$ for the co-hub type knowledge in JEEK via: (1) We initialize $\{W_I^{(i)}, W_S\}$ with $\mathbf{1}_{p \times p}$; (2) $W_{S,j,k} = \frac{1}{\gamma}, \forall j \in NId$ and $k \in \{1, \dots, p\}$ where γ is a hyperparameter. Therefore, the smaller weights for the edge connecting to the node j of all the graphs enforce the co-hub structure.; (3). After this process, each entry of $\{W_I^{(i)}, W_S\}$ equals to either $\frac{1}{\gamma}$ or 1. The below sub-figure of Figure S:4 is an example of the designed W_S .

S:5.3 Case study III: Knowledge of the perturbed hub nodes

Another structure assumption we study is graphs with perturbed nodes. Namely, there exists a set of nodes $NId = \{j | j \in \{1, 2, \dots, p\}\}$ so that there exists i, i'

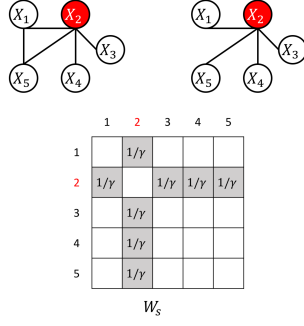


Figure S:1. co-hub. Top: An example of the co-hub node structure. Bottom: The designed W_S for the co-hub structure case.

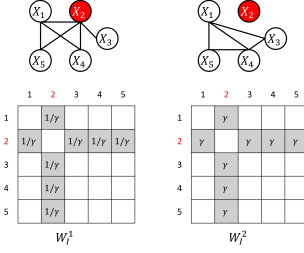


Figure S:2. Perturb hub nodes. Top: An example of the perturbed node structure. Bottom: The designed W_I for the perturbed case.

$\Omega_{j,k}^{(i)} \neq 0$, and $\Omega_{j,k}^{(i')} = 0, \forall k \in \{1, \dots, p\}$. The above sub-figure of Figure S:5 is an example of the perturbed nodes. A so-called JGL-perturb (Mohan et al., 2013) estimator chose $\mathcal{R}'(\cdot) = \sum_{i < i'} P_q((\Omega^{(1)} - \text{diag}(\Omega^{(1)})), \dots, (\Omega^{(K)} - \text{diag}(\Omega^{(K)})))$ in Eq. (5.2). Here $P_q(\cdot)$ has the same definition as mentioned previously. This JGL-perturb formulation also needs a complicated ADMM solution with computationally expensive SVD steps.

To design W_S and $W_I^{(i)}$ for this type of knowledge in JEEK, we use a similar strategy as the above strategy: (1) We initialize $\{W_I^{(i)}, W_S\}$ with $\mathbf{1}_{p \times p}$; we let $W_I^{(i)}_{j,k} = \frac{1}{\gamma}, W_I^{(i')}_{j,k} = \gamma, \forall j \in \text{Id}$ and $k \in 1, \dots, p$. Therefore, the different weights for the edge connecting to the node j in different $W_I^{(i)}$ enforce the node-perturbed structure.; (3). After this process, each entry of $\{W_I^{(i)}, W_S\}$ equals to either $\frac{1}{\gamma}, \gamma$ or 1. The below sub-figure of Figure S:5 is an example of the designed $\{W_I^{(i)}\}$.

S:5.4 Case study IV: Knowledge of group information about nodes

To design W_S and $W_I^{(i)}$ for the group information about a set of nodes, we use a simple three-step strategy: (1) We initialize $\{W_I^{(i)}, W_S\}$ with $\mathbf{1}_{p \times p}$; (2) We let $W_S_{j,k} = \frac{1}{\gamma}, \forall (j, k) \in \text{Id}$ where γ is a hyperparameter. Therefore, the smaller weights for the edge (j, k) in all the graphs favors the edges among nodes in the same group.; (3). After this

process, each entry of $\{W_I^{(i)}, W_S\}$ equals to either $\frac{1}{\gamma}$ or 1. The below sub-figure of Figure S:3 is an example of the designed W_S (extra knowledge is that X_2, X_3, X_4 belong to the same group).

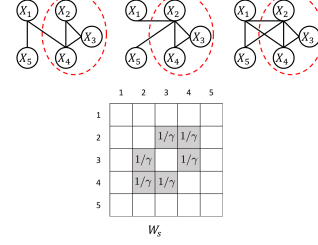


Figure S:3. Co-group

S:6 More about Experimental Setup

S:6.1 Experimental Setup

On four types of datasets, we focus on empirically evaluating JEEK with regard to three aspects: (i) effectiveness, computational speed and scalability in brain connectivity simulation data; (ii) flexibility in incorporating different types of knowledge of known hub nodes in graphs; (iii) effectiveness and computational speed for brain connectivity estimation from real-world fMRI.

S:6.2 Evaluation Metrics

- **AUC-score:** The edge-level false positive rate (FPR) and true positive rate (TPR) are used to measure the difference between the true graphs and the predicted graphs. We obtain FPR vs. TPR curve for each method by tuning over a range of its regularization parameter. We use the area under the FPR -TPR curve (AUC-Score) to compare the predicted versus true graph. Here, $\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$ and $\text{TPR} = \frac{\text{TP}}{\text{FN} + \text{TP}}$. TP (true positive) and TN (true negative) means the number of true edges and non-edges correctly estimated by the predicted network respectively. FP (false positive) and FN (false negative) are the number of incorrectly predicted nonzero entries and zero entries respectively.
- **F1-score:** We first use the edge-level F1-score to compare the predicted versus true graph. Here, $\text{F1} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$, where $\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ and $\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$. The better method achieves a higher F1-score.
- **Time Cost:** We use the execution time (measured in seconds or log(seconds)) for a method as a measure of its scalability. To ensure a fair comparison, we try 30 different λ_n (or λ_2) and measure the total time of execution for

each method. The better method uses less time²

Evaluations: For the first experiment on brain simulation data, we evaluate JEEK and the baseline methods on F1-score and running time cost. For the second experiment, we use AUC-score and running time cost.³ For the third experiment, our evaluation metrics include classification accuracy, likelihood and running time cost.

- The first set of experiments evaluates the speed and scalability of our model JEEK on simulation data imitating brain connectivity. We compare both the estimation performance and computational time of JEEK with the baselines in multiple simulated datasets.
- In the second experiment, we show JEEK’s ability to incorporate knowledge of known hubs in multiple graphs. We also compare the estimation performance and scalability of JEEK with the baselines in multiple simulated datasets.
- Thirdly, we evaluate the ability to import additional knowledge for enhancing graph estimation in a real world dataset. The dataset used in this experiment is a human brain fMRI dataset with two groups of subjects: autism and control. Our choice of this dataset is motivated by recent literature in neuroscience that has suggested many known weights between different regions in human brain as the additional knowledge.

S:6.3 Hyper-parameters:

We need to tune four hyper-parameters v , λ_n , λ_2 and γ :

- v is used for soft-thresholding in JEEK. We choose v from the set $\{0.001i | i = 1, 2, \dots, 1000\}$ and pick a value that makes $T_v(\hat{\Sigma}^{(i)})$ invertible.
- λ_n is the main hyper-parameter that controls the sparsity of the estimated network. Based on our convergence rate analysis in Section 4, $\lambda_n \geq C \sqrt{\frac{\log Kp}{n_{tot}}}$ where $n_{tot} = Kn$ and $n = n_i$. Accordingly, we choose λ_n from a range of $\{0.01 \times \sqrt{\frac{\log Kp}{n_{tot}}} \times i | i \in \{1, 2, 3, \dots, 30\}\}$.
- λ_2 controls the regularization of the second penalty function in JGL-type estimators. We tune λ_2 from the set $\{0.01, 0.05, 0.1\}$ for all experiments and pick the one that gives the best results.

²The machine that we use for experiments is an AMD 64-core CPU with a 256GB memory.

³We cannot use AUC-score for the first set of experiments as the baseline NAK only gives us the best adjacency matrix after tuning over their hyperparameters. It does not provide an option for tuning the λ_n .

- γ is a hyperparameter used to design the $W_I^{(i)}$, W_S (5). The value of γ intuitively indicates the confidence of the additional knowledge weights. In the second experiment, we choose $\gamma = \{2, 4, 10\}$.

S:7 More about Experimental Results

S:7.1 More Experiment: Simulate Samples with Known Hubs as Knowledge

In this set of experiments, we show empirically JEEK’s ability to model knowledge of known hub nodes across multiple sGGMs and its advantages in scalability and effectiveness. We generate multiple simulated Gaussian datasets for both the co-hub and perturbed-hub graph structures.

Simulation Protocol to generate simulated datasets: We generate multiple sets of synthetic multivariate-Gaussian datasets. First, we generate random graphs following the Random Graph Model (Rothman et al., 2008). This model assumes $\Omega^{(i)} = \mathbf{B}_I^{(i)} + \mathbf{B}_S + \delta I$, where each off-diagonal entry in $\mathbf{B}^{(i)}$ is generated independently and equals 0.5 with probability $0.1i$ and 0 with probability $1 - 0.1i$. The shared part \mathbf{B}_S is generated independently and equal to 0.5 with probability 0.1 and 0 with probability 0.9. δ is selected large enough to guarantee positive definiteness. We generate cohub and perturbed structure simulations, using the following data generation models:

- **Random Graphs with cohub nodes:** After we generate the random graphs using the aforementioned Random Graph Model, we randomly generate a set of nodes $NId = \{j | j \in \{1, 2, \dots, p\}\}$ as the cohub nodes among all the random graphs. The cardinal number of this set equals to $5\%p$. For each of these nodes j , we randomly select 90% edges $E_j = \{(j, k) | k \in \{1, 2, \dots, p\}\}$ to be included in the graph. Then we set $\Omega_{j,k}^{(i)} = \Omega_{k,j}^{(i)} = 0.5, \forall i \in \{1, 2, \dots, K\}$ and $(j, k) \in E_j$.
- **Random Graphs with perturbed nodes:** After we generate the random graphs using the aforementioned Random Graph Model, we randomly generate a set of nodes $NId = \{j | j \in \{1, 2, \dots, p\}\}$ as the perturbed hub nodes for the random graphs. The cardinal number of this set equals to $5\%p$. For all graphs $\{\Omega^{(i)} | i \text{ is odd}\}$, for each of these nodes $j \in NId$, we randomly select 90% edges $E_j = \{(j, k) | k \in \{1, 2, \dots, p\}\}$ to be included in the graph. We set $\Omega_{j,k}^{(i)} = \Omega_{k,j}^{(i)} = 0.5, \forall \text{ odd } i \in \{1, 2, \dots, K\}$ and $(j, k) \in E_j$. For all graphs $\{\Omega^{(i)} | i \text{ is even}\}$ and nodes $j \in NId$, we randomly select 10% edges $E'_j = \{(j, k) | k \in \{1, 2, \dots, p\}\}$ to be included in the graph. We set $\Omega_{j,k}^{(i)} = \Omega_{k,j}^{(i)} = 0.5, \forall \text{ even } i \in \{1, 2, \dots, K\}$ and $(j, k) \in E'_j$. This creates a perturbed node structure in the multiple graphs.

Experimental baselines: We employ JGL-node for cohub and perturbed hub node structure (JGL-hub and JGL-perturb respectively) and W-SIMULE as the baselines for this set of experiments. The weights in $\{W_I^{tot}, W_S^{tot}\}$ are designed by the strategy mentioned in Section S:5.

Experiment Results: We assess the performance of JEEK in terms of effectiveness (AUC score) and scalability (computational time cost) through baseline comparison as follows:

(a) Effectiveness: We plot the AUC-score for a number of multiple simulated datasets generated by varying the number of features p , the number of tasks K and the number of samples n . We calculate AUC by varying λ_n . For the JGL estimator, we additionally vary λ_2 and select the best AUC (section S:6.1). In Figure S:4 (a) and Figure S:4 (b), we plot the AUC-Score for the cohub node structure vs varying p and K , respectively. Figure S:5 (a) and Figure S:5 (b) plot the same for the perturbed node structure. In Figure S:4 (a) and Figure S:5 (a), we vary p in the set $\{100, 200, 300, 400, 500\}$ and set $K = 2$ and $n = p/2$. For $p > 300$ and $n = p/2$, W-SIMULE takes more than one month and JGL takes more than one day. Therefore we can not show their results for $p > 300$. For both the cohub and perturbed node structures, JEEK consistently achieves better AUC-score than the baseline methods as p is increased. For Figure S:4(b) and Figure S:5 (b), we vary K in the set $\{2, 3, 4\}$ and set $p = 200$ and $n = p/2$. JEEK consistently has a higher AUC-score than the baselines JGL and W-SIMULE as K is increased.

(b) Scalability: In Figure S:4 (c) and (d), we plot the computational time cost for the cohub node structure vs the number of features p and the number of tasks K , respectively. Figure S:5 (c) and (d) plot the same for the perturbed node structure. We interpolate the points of computation time of each estimator into curves. For each simulation case, the computation time for each estimator is the summation of a method’s execution time over all values of λ_n . In Figure S:4(c) and Figure S:5(c), we vary p in the set $\{100, 200, 300, 400, 500\}$ and set $K = 2$ and $n = p/2$. When $p > 300$ and $n = p/2$, W-SIMULE takes more than one month and JGL takes more than one day. Hence, we have omitted their results for $p > 300$. For both the cohub and perturbed node structures, JEEK is consistently more than 5 times faster as p is increased. In Figure S:4(d) and Figure S:5 (d), we vary K in the set $\{2, 3, 4\}$ and fix $p = 200$ and $n = p/2$. JEEK is 50 times faster than the baselines for all cases with $p = 200$ and as K is increased. In summary, JEEK is on an average more than 10 times faster than all the baselines.

(c) Stability of Results when varying W matrices: Additionally, to account for JEEK’s explicit structure assumption, we also vary the ratio of known hub nodes to the total

number of hub nodes. The known hub nodes are used to design the $\{W_I^i, W_S\}$ matrices(details in Section 5). In Figure S:6(a) and (b), AUC for JEEK increases as the ratio of the number of known to total hub nodes increases. The initial increase in AUC is particularly significant as it confirms that JEEK is effective in harvesting additional knowledge for multiple sGGMs. The increase in AUC is particularly significant in the perturbed node case (Figure S:6(b)). The AUC for the hub case does not have a correspondingly large increase with an increase in ratio because the total number of hub nodes are only 5% of the total nodes. In comparison, an increase in this ratio leads to a more significant increase in AUC because the perturbed node assumption has more information than the cohub node structure. We show in Figure S:6(c) and (d) that the computational cost is largely unaffected by this ratio for both the cohub and perturbed node structure.

We also empirically check how the parameter r in the designed knowledge weight matrices influences the performance. In Figure S:7(a) and (b), we show that the designed strategy for including additional knowledge as W is not affected by variations of γ . We vary γ in the set of $\{2, 4, 10\}$. In summary, the AUC-score(Figure S:7(a),(b)) and computational time cost(Figure S:7(c),(d)) remains relatively unaffected by the changes in γ for both co-hub and perturbed-hub case.

S:7.2 More Experiment: Gene Interaction Network from Real-World Genomics Data

Next, we apply JEEK and the baselines on one real-world biomedical data: gene expression profiles describing many human samples across multiple cancer types aggregated by (McCall et al., 2011).

Advancements in genome-wide monitoring have resulted in enormous amounts of data across most of the common cell contexts, like multiple common cancer types (Network et al., 2011). Complex diseases such as cancer are the result of multiple genetic and epigenetic factors. Thus, recent research has shifted towards the identification of multiple genes/proteins that interact directly or indirectly in contributing to certain disease(s). Structure learning of sGGMs on such heterogeneous datasets can uncover statistical dependencies among genes and understand how such dependencies vary from normal to abnormal or across different diseases. These structural variations are highly likely to be contributing markers that influence or cause the diseases.

Two major cell contexts are selected from the human expression dataset provided by (McCall et al., 2011): leukemia cells (including 895 samples and normal blood cells (including 227 samples)). Then we choose the top 1000 features from the total 12,704 features (ranked by variance) and perform graph estimation on this two-task dataset. We explore two type of knowledge in the experiments.

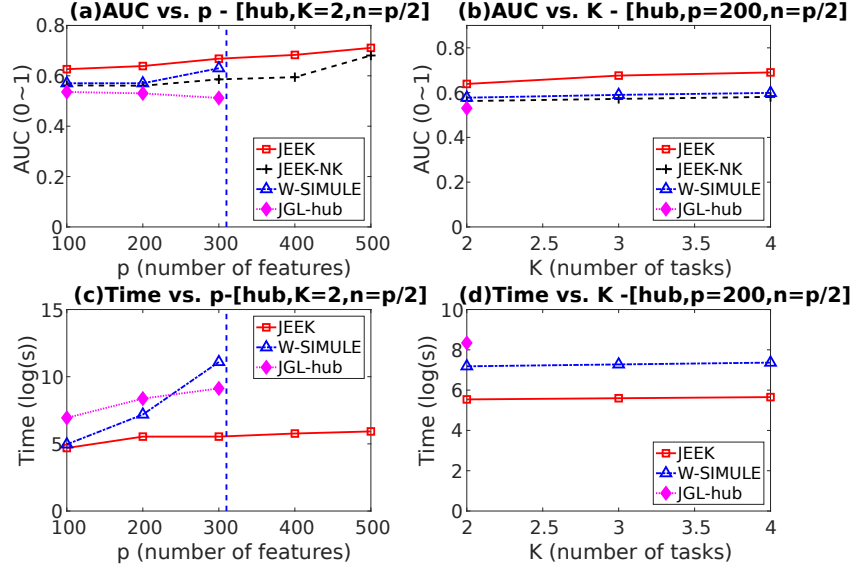


Figure S:4. Cohub node structure: (a) AUC-score vs the number of features (p). (b) AUC-score vs the number of tasks (K). (c) Time cost (log(seconds)) vs the number of features (p). (d) Time cost (log(seconds)) vs the number of tasks (K). For $p > 300$ and $n = p/2$ W-SIMULE takes more than one month and JGL takes more than one day (indicated by dotted blue line). JGL package can only run for $K = 2$.

The first kind (DAVID) is about the known group information about nodes, such as genes belonging to the same biological pathway or cellular location. We use the popular “functional enrichment” analysis tool DAVID (Da Wei Huang & Lempicki, 2008) to get a set of group information about the 1000 genes. Multiple different types of groups are provided by DAVID and we pick the co-pathway. We only use the grouping information covering 20% of the nodes (randomly picked from 1000). The derived dependency graphs are compared by using the number of predicted edges being validated by three major existing protein/gene interaction databases (Prasad et al., 2009; Orchard et al., 2013; Stark et al., 2006) (average over both cell contexts).

The second type (PPI) is using existing known edges as the knowledge, like the known protein interaction databases for discovering gene networks (a semi-supervised setting for such estimations). We use three major existing protein/gene interaction databases (Prasad et al., 2009; Orchard et al., 2013; Stark et al., 2006). We only use the known interaction edge information covering 20% of the nodes (randomly picked from 1000). The derived dependency graphs are compared by using the number of predicted edges that are not part of the known knowledge and are being validated by three major existing protein/gene interaction databases (Prasad et al., 2009; Orchard et al., 2013; Stark et al., 2006) (average over both cell contexts).

We would like to point out that the interactions JEEK and baselines find represent statistical dependencies between

genes that vary across multiple cell types. There exist many possibilities for such interactions, including like physical protein-protein interactions, regulatory gene pairs or signaling relationships. Therefore, we combine multiple existing databases for a joint validation. The numbers of matches between interactions in databases and those edges predicted by each method have been shown as the y -axis in Figure S:8(c). It clearly shows that JEEK consistently outperforms two baselines.

S:7.3 More Experiment: Simulated Samples about Brain Connectivity with Distance as Knowledge

In this set of experiments, we confirm JEEK’s ability to harvest additional knowledge using brain connectivity simulation data. Following (Bu & Lederer, 2017), we employ the known Euclidean distance between brain regions as additional knowledge W to generate simulated datasets. To generate the simulated graphs, we use $p_{j,k} = \text{inv.logit}(10 - W_{j,k}/3)$ as the probability of an edge between nodes j and k in the graphs, where $W_{j,k}$ is the Euclidean distance between regions j and k of the brain.

The generate datasets all have $p = 116$ corresponding to the number of brain regions in the distance matrix shared by (Bu & Lederer, 2017). We vary K from the set $\{2, 3, 4\}$ with $n = p/2$. The F1-scores for JEEK, JEEK-NK and W-SIMULE is the best F1-score after tuning over λ_n . The hyperparameter tuning for NAK is done by the package itself.

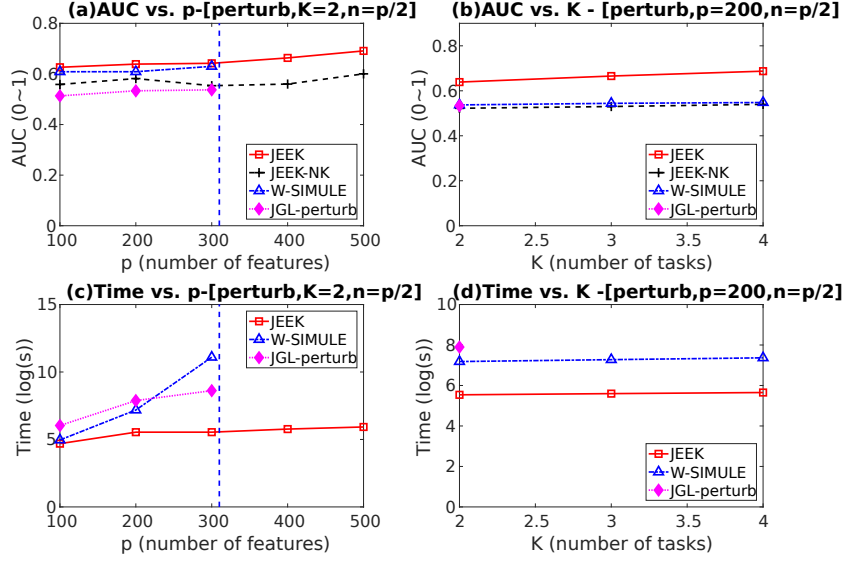


Figure S:5. Perturbed node structure: (a) AUC-score vs the number of features (p). (b) AUC-score vs the number of tasks (K). (c) Time cost (log(seconds)) of JEEK and the baseline methods vs the number of features (p). (d) Time cost (log(seconds)) vs the number of tasks (K). for $p > 300$ and $n = p/2$, W-SIMULE takes more than one month and JGL takes more than one day (indicated by dotted blue line). JGL package can only run for $K = 2$.

Simulated brain data generation model: We generate multiple sets of synthetic multivariate-Gaussian datasets. To imitate brain connectivity, we use the Euclidean distance between the brain regions as additional knowledge W where $W_{j,k}$ is the Euclidean distance between regions j and k . We fix $p = 116$ corresponding to the number of brain regions (Bu & Lederer, 2017). We generate the graph $\Omega^{(i)}$ following $\Omega^{(i)} = \mathbf{B}_I^{(i)} + \mathbf{B}_S + \delta I$, where each off-diagonal entry in $\mathbf{B}_I^{(i)}$ is generated independently and equals 0.5 with probability $p_{j,k} = \text{inv.logit}(10 - W_{j,k}/3)$ and 0 with probability $1 - p_{j,k}$ (Bu & Lederer, 2017). Similarly, the shared part \mathbf{B}_S is generated independently and equal to 0.5 with probability $p_{j,k} = \text{inv.logit}(10 - W_{j,k}/3)$ and 0 with probability $1 - p_{j,k}$. δ is selected large enough to guarantee the positive definiteness. This choice ensures there are more direct connections between close regions, effectively simulating brain connectivity. For each case of simulated data generation, we generate K blocks of data samples following the distribution $N(0, (\Omega^{(i)})^{-1})$. Details see Section S:6.1.

Experimental baselines: We choose W-SIMULE, NAK and JEEK with no additional knowledge (JEEK-NK) as the baselines. (see Section 5).

Experiment Results: We compare JEEK with the baselines regarding two aspects— (a) Scalability (Computational time cost), and (b) Effectiveness (F1-score). Figure S:9(a) and Figure S:9(b) respectively show the F1-score vs. computational time cost with varying number of tasks K and the number of samples n . In these experiments, $p = 116$

corresponding to the number of brain regions in the distance matrix provided by (Bu & Lederer, 2017). In Figure S:9(a), we vary K in the set $\{2, 3, 4\}$ with $n = p/2$. In Figure S:9(b), we vary n in the set $\{p/2, p, 2p\}$ and fix $K = 2$. The F1-score plotted for JEEK, JEEK-NK and W-SIMULE is the best F1-score after tuning over λ_n . The hyperparameter tuning for NAK is done by the package itself. For each simulation case, the computation time for each estimator is the summation of a method’s execution time over all values of λ_n . The points in the top left region of Figure S:9 indicate higher F1-score and lower computational cost. Clearly, JEEK outperforms its baselines as all JEEK points are in the top left region of Figure S:9. JEEK has a consistently higher F1-Score and is almost 6 times faster than W-SIMULE in the high dimensional case. JEEK performs better than JEEK-NK, confirming the advantage of integrating additional knowledge in graph estimation. While NAK is fast, its F1-Score is nearly 0 and hence, not useful for multi-sGGM estimation.

S:7.4 More Experiment: Brain Connectivity Estimation from Real-World fMRI

Experimental Baselines: We choose W-SIMULE as the the baseline in this experiment. We also compare JEEK to JEEK-NK and W-SIMULE-NK to demonstrate the need for additional knowledge in graph estimation.

ABIDE Dataset: This data is from the Autism Brain Imaging Data Exchange (ABIDE) (Di Martino et al., 2014), a publicly available resting-state fMRI dataset. The ABIDE

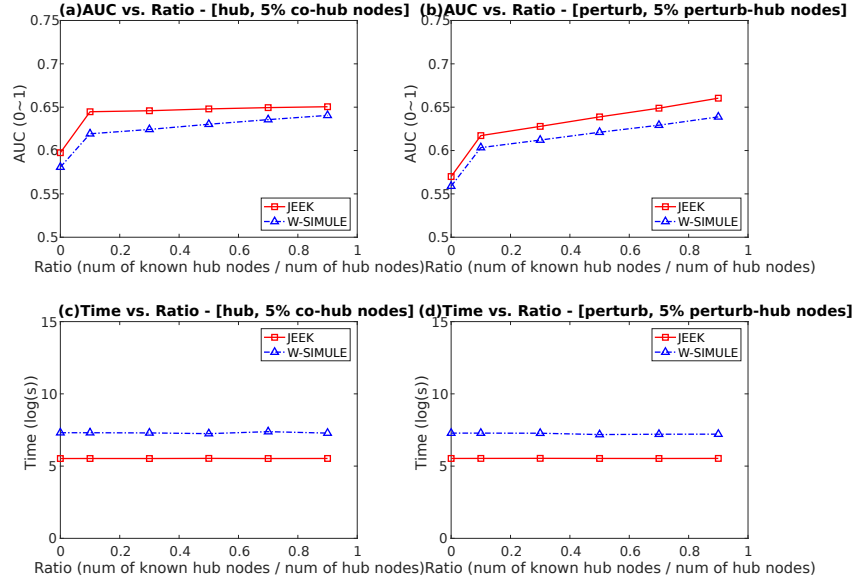


Figure S:6. AUC-Score vs. ratio of number of known hub nodes to number of total hub nodes for (a) Cohub node structure (b) perturbed node structure. Computational Time Cost vs. ratio of number of known hub nodes to number of total hub nodes for (a) Cohub node structure (b) perturbed node structure.

data aims to understand human brain connectivity and how it reflects neural disorders (Van Essen et al., 2013). The data is retrieved from the Preprocessed Connectomes Project (Craddock, 2014), where preprocessing is performed using the Configurable Pipeline for the Analysis of Connectomes (CPAC) (Craddock et al., 2013) without global signal correction or band-pass filtering. After preprocessing with this pipeline, 871 individuals remain (468 diagnosed with autism). Signals for the 160 (number of features $p = 160$) regions of interest (ROIs) in the often-used Dosenbach Atlas (Dosenbach et al., 2010) are examined.

Distance as Additional Knowledge: To select the weights $\{W_I^{(i)}, W_S\}$, two separate spatial distance matrices W were derived from the Dosenbach atlas. The first, referred to as *anatomical* ^{i} , gives each ROI one of 40 well-known, anatomic labels (e.g. “basal ganglia”, “thalamus”). Weights $W_{j,k}$ take the low value i if two ROIs have the same label, and the high value $10 - i$ otherwise. The second additional knowledge matrix, referred to as *dist* ^{i} , sets the weight of each edge ($W_{j,k}$) to its spatial length, in MNI space⁴, raised to the power i . Then $W_I^{(i)} = W_S = W$.

Cross-validation: Classification is performed using the 3-fold cross-validation suggested by the literature (Poldrack et al., 2008)(Varoquaux et al., 2010). The subjects are randomly partitioned into three equal sets: a training set,

⁴MNI space is a coordinate system used to refer to analogous points on different brains.

a validation set, and a test set. Each estimator produces $\hat{\Omega}^{(1)} - \hat{\Omega}^{(2)}$ using the training set. Then, these differential networks are used as inputs to linear discriminant analysis (LDA), which is tuned via cross-validation on the validation set. Finally, accuracy is calculated by running LDA on the test set. This classification process aims to assess the ability of an estimator to learn the differential patterns of the connectome structures. We cannot use NAK to perform classification for this task, as NAK outputs only an adjacency matrix, which cannot be used for estimation using LDA.

Parameter variation: The results are fairly robust to variations of the W . (see Table S:1). The effect of changing W seems to have a fairly small effect on the log-likelihood of the model. This is likely because both penalize picking physically long edges, which agrees with observations from neuroscience. The *dist* W effectively encourages the selection of short edges, and the *anatomical* W also has substantial spatial localization.

References

- Bickel, P. J. and Levina, E. Covariance regularization by thresholding. *The Annals of Statistics*, pp. 2577–2604, 2008.
- Bu, Y. and Lederer, J. Integrating additional knowledge into estimation of graphical models. *arXiv preprint arXiv:1704.02739*, 2017.
- Craddock, C. Preprocessed connectomes project: open shar-

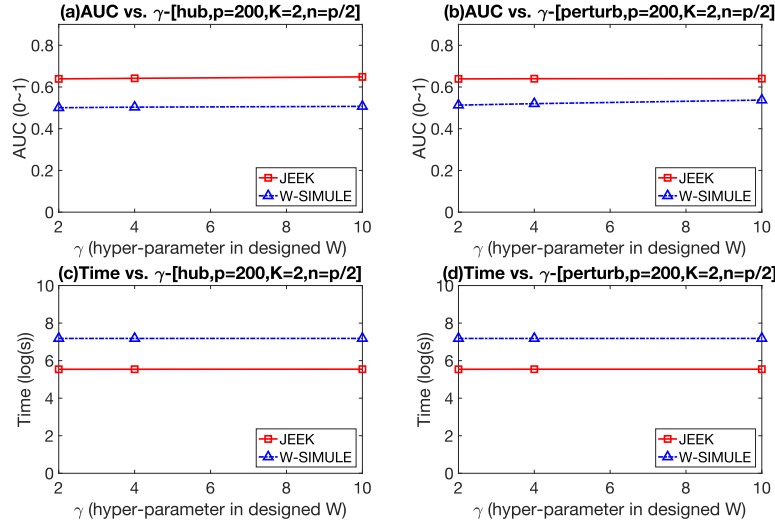


Figure S:7. AUC-Score vs. γ (a) Cohub node structure for (b) perturbed node structure. Computational Time Cost vs. γ for (a) Cohub node structure (b) perturbed node structure.

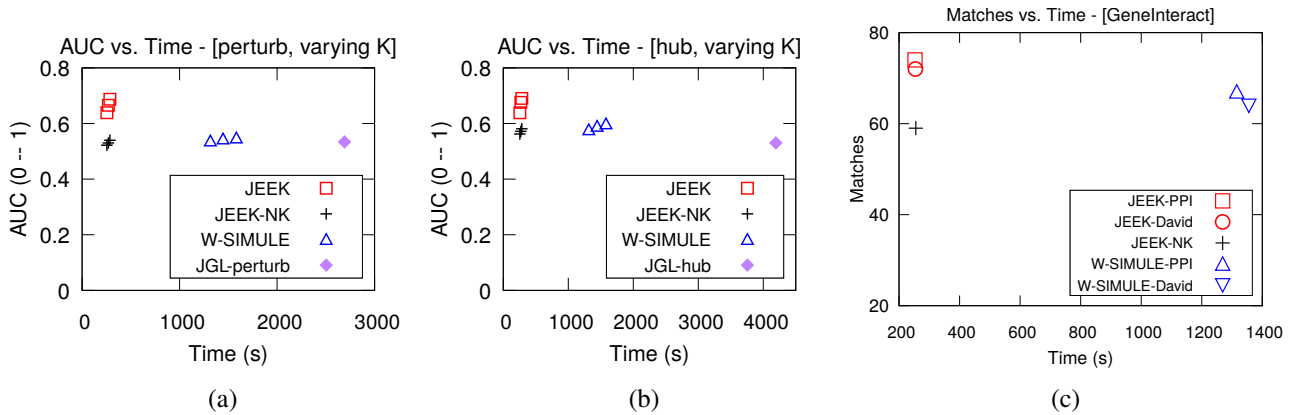


Figure S:8. (a)(b) Performance comparison on simulation Datasets about hubs: AUC vs. Time when varying number of tasks K . (a) is the perturbed hub cases and (b) is for the co-hub cases. (c) Performance comparison on one real-world gene expression dataset with two cell types. Two type knowledge are used to cover one fifth of the nodes, therefore each method corresponds to two performance points.

ing of preprocessed neuroimaging data and derivatives. In *61st Annual Meeting. AACAP*, 2014.

Craddock, C., Sikka, S., Cheung, B., Khanuja, R., Ghosh, S., Yan, C., Li, Q., Lurie, D., Vogelstein, J., Burns, R., et al. Towards automated analysis of connectomes: The configurable pipeline for the analysis of connectomes (c-pac). *Front Neuroinform*, 42, 2013.

Da Wei Huang, B. T. S. and Lempicki, R. A. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nature protocols*, 4(1):44–57, 2008.

Danaher, P., Wang, P., and Witten, D. M. The joint graphical lasso for inverse covariance estimation across multiple

classes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2013.

Di Martino, A., Yan, C.-G., Li, Q., Denio, E., Castellanos, F. X., Alaerts, K., Anderson, J. S., Assaf, M., Bookheimer, S. Y., Dapretto, M., et al. The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism. *Molecular psychiatry*, 19(6): 659–667, 2014.

Dosenbach, N. U., Nardos, B., Cohen, A. L., Fair, D. A., Power, J. D., Church, J. A., Nelson, S. M., Wig, G. S., Vogel, A. C., Lessov-Schlaggar, C. N., et al. Prediction of individual brain maturity using fmri. *Science*, 329(5997): 1358–1361, 2010.

McCall, M. N., Uppal, K., Jaffee, H. A., Zilliox, M. J., and

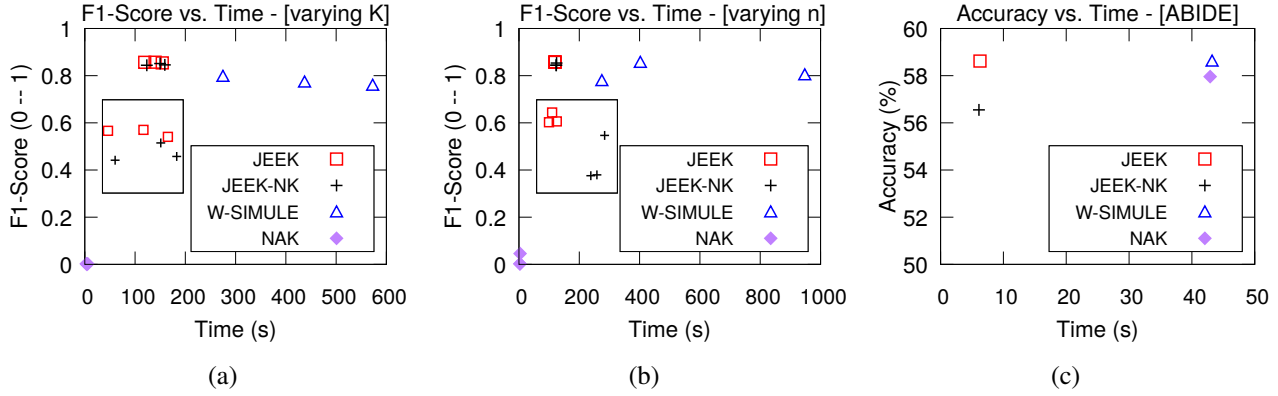


Figure S:9. Experimental Results on Simulated Brain Datasets and on ABIDE. (a) Performance obtained on simulated brain samples with respect to F1-score vs. computational time cost when varying the number of tasks K . (b) Performance obtained on simulated brain samples with respect to F1-score vs. computational time cost when varying the number of samples n . In both (a) and (b) the smaller box shows an enlarged view comparing JEEK and JEEK-NK points. All JEEK points are in the top left region indicating higher F1-score and lower computational cost. (c). On ABIDE, JEEK outperforms the baseline methods in both classification accuracy and running time cost. JEEK and JEEK-NK points in the top left region and JEEK points are higher in terms of y -axis positions.

Table S:1. Variations of the W and multi-task component yield fairly stable results.

Prior	Sparsity=8%		Sparsity=16%	
	Log-Likelihood	Test Accuracy	Log-Likelihood	Test Accuracy
No Additional Knowledge	-294.34	0.56	-283.27	0.55
$dist$	-289.12	0.53	-285.69	0.55
$dist^2$	-283.78	0.54	-282.92	0.54
$anatomical^1$	-292.42	0.56	-289.34	0.57
$anatomical^2$	-291.29	0.58	-285.63	0.56

Irizarry, R. A. The gene expression barcode: leveraging public data repositories to begin cataloging the human and murine transcriptomes. *Nucleic acids research*, 39 (suppl 1):D1011–D1015, 2011.

Mohan, K., London, P., Fazel, M., Lee, S.-I., and Witten, D. Node-based learning of multiple gaussian graphical models. *arXiv preprint arXiv:1303.5145*, 2013.

Negahban, S., Yu, B., Wainwright, M. J., and Ravikumar, P. K. A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. In *Advances in Neural Information Processing Systems*, pp. 1348–1356, 2009.

Network, C. G. A. R. et al. Integrated genomic analyses of ovarian carcinoma. *Nature*, 474(7353):609–615, 2011.

Orchard, S., Ammari, M., Aranda, B., Breuza, L., Briganti, L., Broackes-Carter, F., Campbell, N. H., Chavali, G., Chen, C., Del-Toro, N., et al. The MIntAct project IntAct as a common curation platform for 11 molecular interaction databases. *Nucleic acids research*, pp. gkt1115, 2013.

Park, T. and Casella, G. The bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.

Poldrack, R. A., Fletcher, P. C., Henson, R. N., Worsley, K. J., Brett, M., and Nichols, T. E. Guidelines for reporting an fmri study. *Neuroimage*, 40(2):409–414, 2008.

Prasad, T. K., Goel, R., Kandasamy, K., Keerthikumar, S., Kumar, S., Mathivanan, S., Telikicherla, D., Raju, R., Shafreen, B., Venugopal, A., et al. Human protein reference database?2009 update. *Nucleic acids research*, 37 (suppl 1):D767–D772, 2009.

Ravikumar, P., Wainwright, M. J., Raskutti, G., Yu, B., et al. High-dimensional covariance estimation by minimizing l1-penalized log-determinant divergence. *Electronic Journal of Statistics*, 5:935–980, 2011.

Rothman, A. J., Bickel, P. J., Levina, E., Zhu, J., et al. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494–515, 2008.

Rothman, A. J., Levina, E., and Zhu, J. Generalized thresholding of large covariance matrices. *Journal of the American Statistical Association*, 104(485):177–186, 2009.

Stark, C., Breitkreutz, B.-J., Reguly, T., Boucher, L., Breitkreutz, A., and Tyers, M. Biogrid: a general repository for interaction datasets. *Nucleic acids research*, 34 (suppl.1):D535–D539, 2006.

- Van Essen, D. C., Smith, S. M., Barch, D. M., Behrens, T. E., Yacoub, E., Ugurbil, K., Consortium, W.-M. H., et al. The wu-minn human connectome project: an overview. *Neuroimage*, 80:62–79, 2013.
- Varoquaux, G., Gramfort, A., Poline, J.-B., and Thirion, B. Brain covariance selection: better individual functional connectivity models using population prior. In *Advances in neural information processing systems*, pp. 2334–2342, 2010.
- Wainwright, M. J. and Jordan, M. I. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1-2):1–305, 2008.
- Wang, B., Gao, J., and Qi, Y. A fast and scalable joint estimator for learning multiple related sparse gaussian graphical models. In *Artificial Intelligence and Statistics*, pp. 1168–1177, 2017a.
- Wang, B., Singh, R., and Qi, Y. A constrained l1 minimization approach for estimating multiple sparse gaussian or nonparanormal graphical models. *Machine Learning*, 106(9-10):1381–1417, 2017b.
- Watts, D. J. and Strogatz, S. H. Collective dynamics of ‘small-world’ networks. *nature*, 393(6684):440–442, 1998.
- Yang, E., Lozano, A. C., and Ravikumar, P. Elementary estimators for graphical models. In *Advances in Neural Information Processing Systems*, pp. 2159–2167, 2014.