

Paired data

Paired data problem: every variable is uniquely associated with a homologous, or twin, variable.

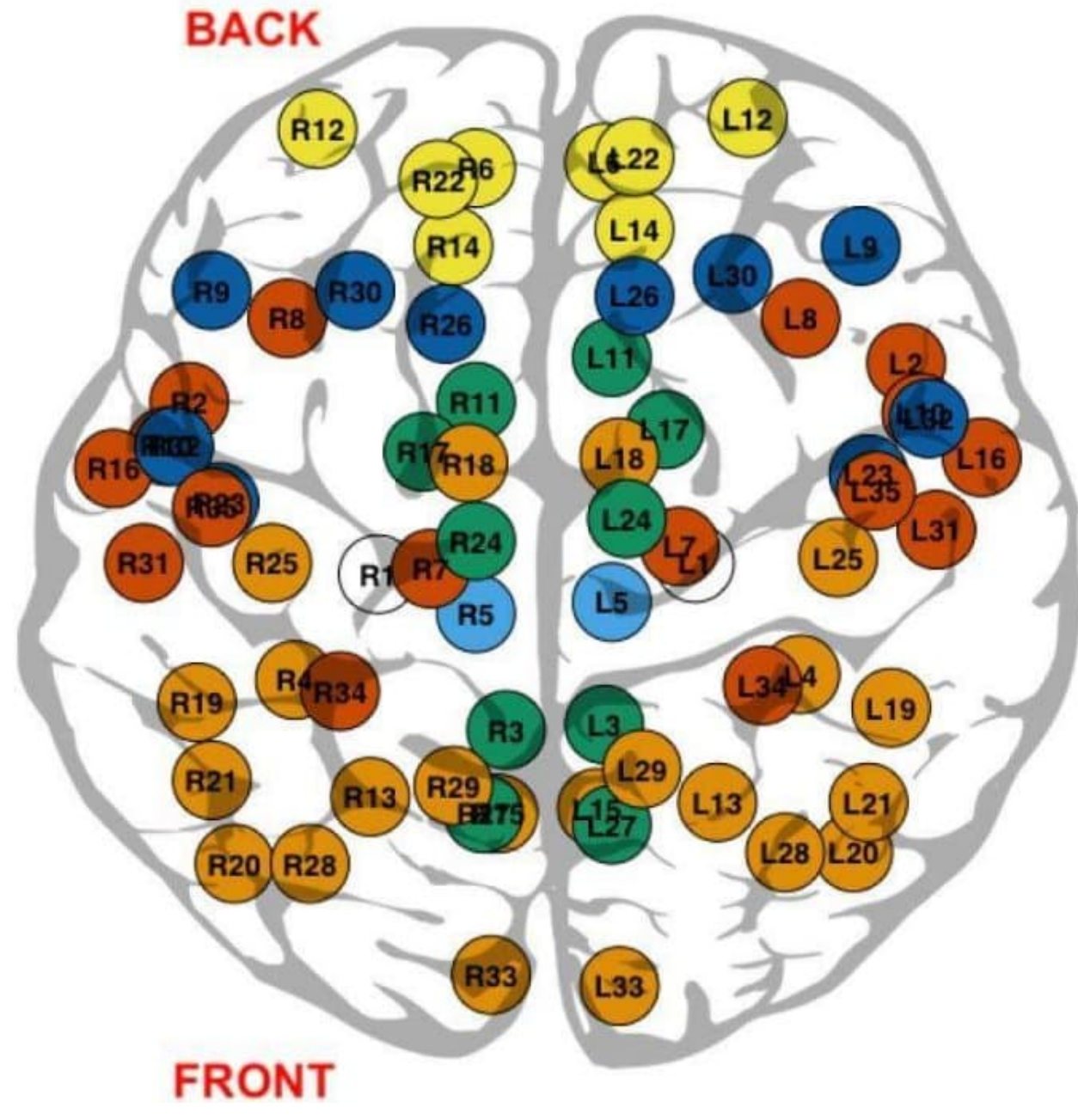


Figure 1. Example of ROI locations on the brain. Every ROI on the left hemisphere is associated with an ROI on the right hemisphere, which gives the pairs $(L_i, R_i)_{i=1,\dots,35}$. Different colors correspond to distinct brain regions.

Hence, for paired data, \mathbf{Y}_V can be partitioned as $(\mathbf{Y}_L, \mathbf{Y}_R)^T$, and we consider and assume that $L = \{1, \dots, q\}$ and $R = \{1', \dots, q'\}$ where $i' = q + i$ and $q = p/2$ so that Y_i is homologous to $Y_{i'}$ with $1 \leq i \leq q$.

Gaussian graphical models (GGMs)

Let $G = (V, E)$ be an undirected graph with the vertex set V and the edge set E . Then, \mathbf{Y}_V is said to satisfy the Gaussian graphical model if $\mathbf{Y}_V \sim \mathcal{N}(\mu, \Sigma)$ and \mathbf{Y}_V is Markov w.r.t G , that is $(i, j) \notin E$ implies $\theta_{ij} = 0$ where $\Theta = (\theta_{ij})_{i,j \in V} = \Sigma^{-1}$.

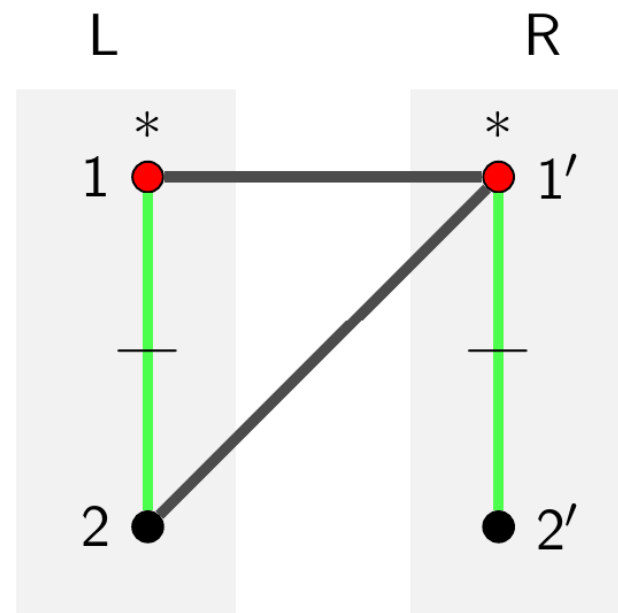
Coloured GGMs for paired data

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a coloured version of G where \mathcal{V} is a partition of V into vertex colour classes; similarly, \mathcal{E} is a partition of E into edge colour classes.

Coloured graphs for paired data (PD-CGs)

The PD-CG \mathcal{G} contains two types of color classes:

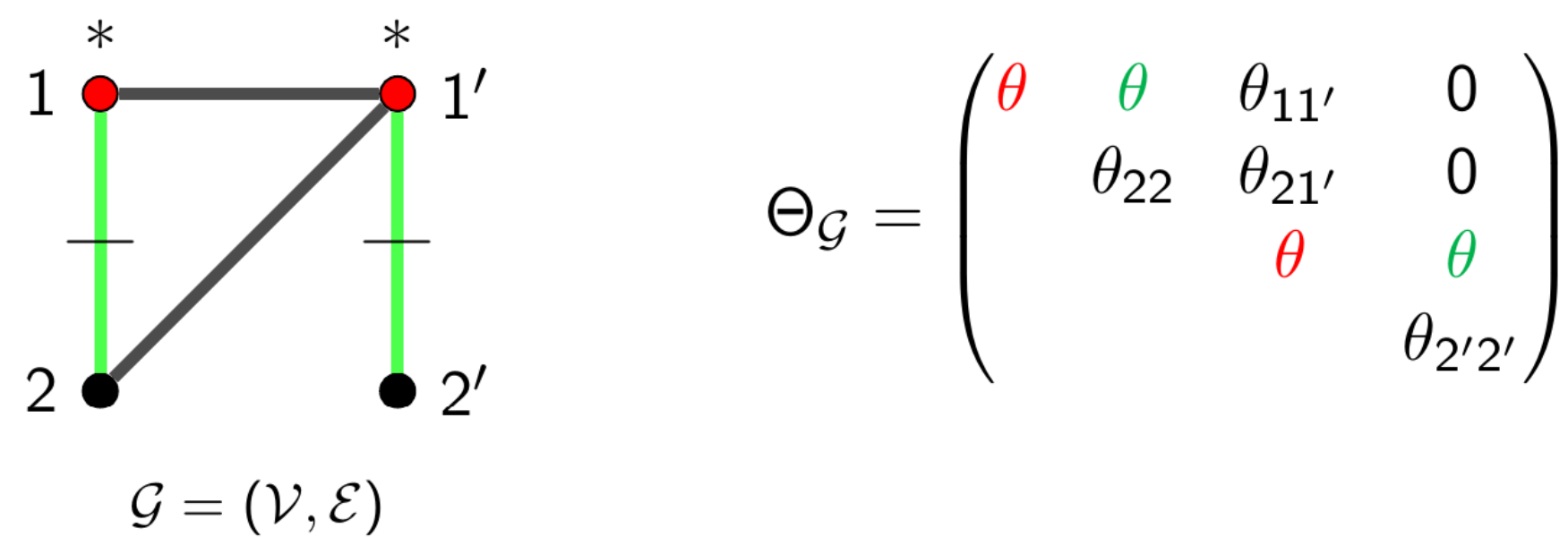
- **atomic class** that is a color class of cardinality one;
- **twin-pairing class** that is a color class containing a pair of twin vertices or a pair of twin edges.



Example. Consider the PD-CG $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with
 $\mathcal{V} = \{ \underbrace{\{1, 1'\}}_{\text{twin-pairing}}, \underbrace{\{2\}, \{2'\}}_{\text{atomic}} \}$, $\mathcal{E} = \{ \underbrace{\{(1, 2), (1', 2')\}}_{\text{twin-pairing}}, \underbrace{\{(1, 1')\}, \{(2', 1')\}}_{\text{atomic}} \}$.

RCON models for paired data (PD-RCONs)

PD-RCON models are Gaussian graphical models with additional equality constraints on the concentration matrix implied by a PD-CG $\mathcal{G} = (\mathcal{V}, \mathcal{E})$.



Challenges

Learn the graphical models for paired data:

1. learn the structure of the graph;
2. learn the colourings of the vertices;
3. learn the colourings of the edges both between and across the left and the right parts of the network.

Difficulties

1. Dimension of the search space, for example:

$$\begin{matrix} 1 \\ \text{complete graph} \\ \text{on } p \text{ vertices} \end{matrix} \ll \begin{matrix} 2^{(p/2)^2} \\ \text{complete graphs} \\ \text{for paired data} \end{matrix}$$

2. The exploration of the space:

- the structure of the search space forms a lattice but it is not distributive,
- the neighbors of a model cannot be efficiently specified.

Structure of models space of PD-CGMs

It is useful to embed search spaces with a partial order. Naturally, the order is the **model inclusion order**: a model is “larger” than any of its submodels.

Consider two PD-CGMs characterized by $\mathcal{G} = (\mathcal{V}_G, \mathcal{E}_G)$ and $\mathcal{H} = (\mathcal{V}_H, \mathcal{E}_H)$. Then, following [2], $\mathcal{G} \preceq_s \mathcal{H}$ if and only if

- $E_H \supseteq E_G$,
 - $\mathcal{V}_H \preceq_f \mathcal{V}_G$,
 - $\mathcal{E}_H \preceq_f \mathcal{E}_G \cup \{\{E_H \setminus E_G\}\}$,
- where \preceq_f is the *refinement order* and E_G, E_H are the sets of uncolored edges of \mathcal{G}, \mathcal{H} , respectively.

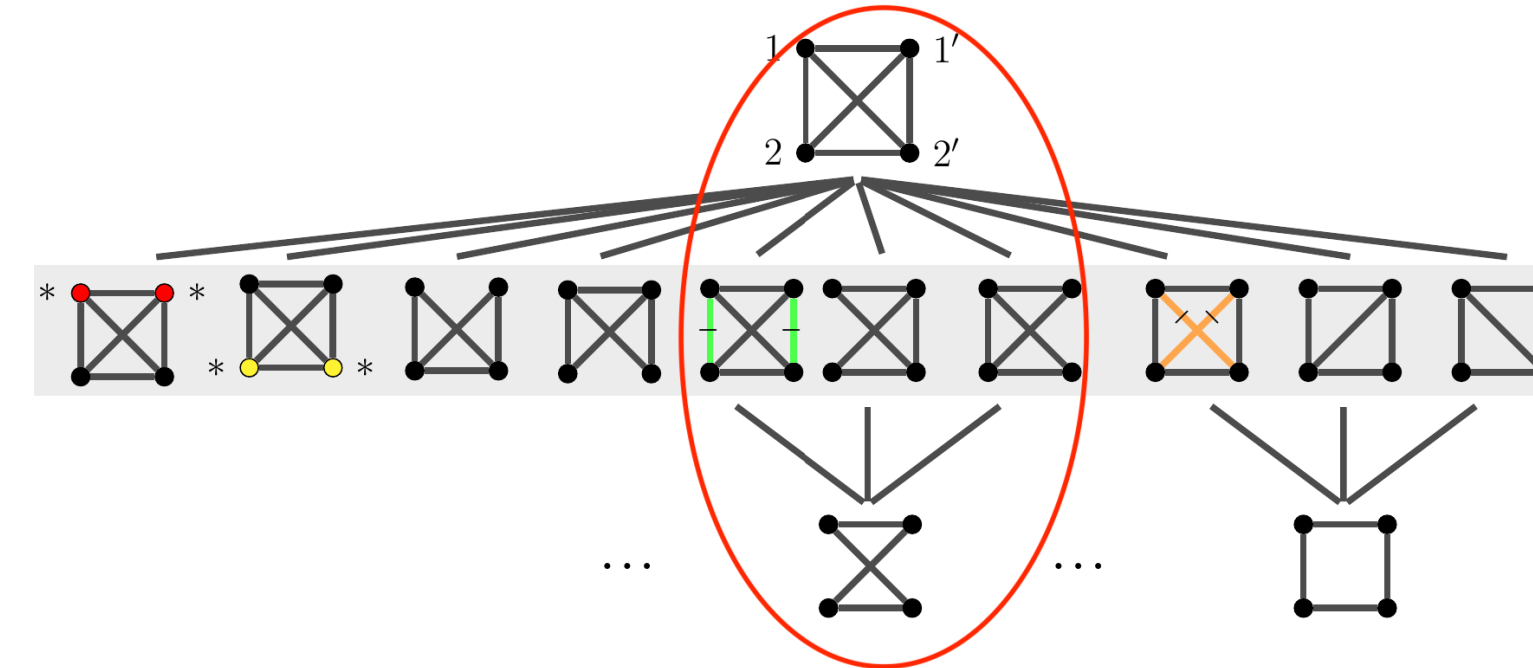


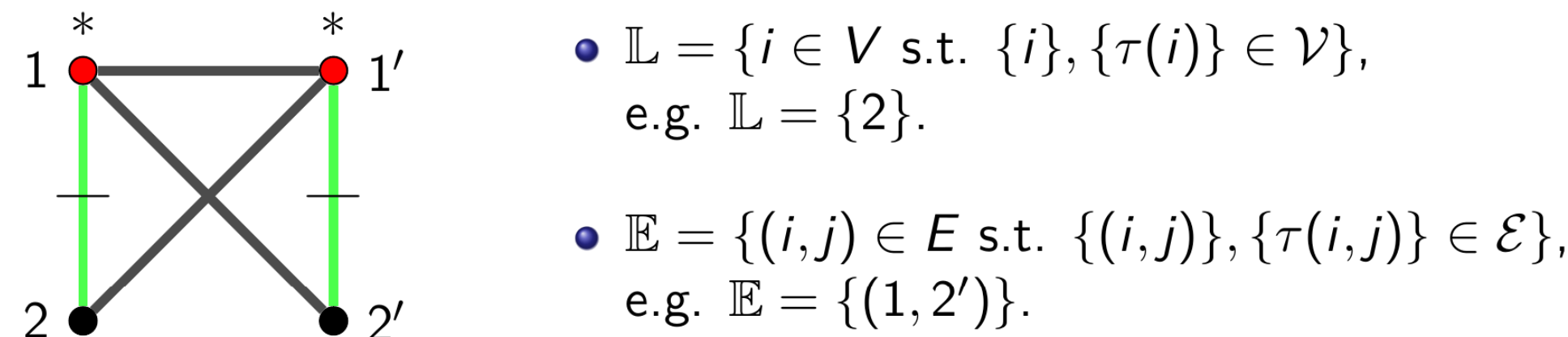
Figure 2. A part of Hasse diagram of the lattice structure of PD-CGs with 4 vertices based on the model inclusion order. The highlighted graphs are the neighbors of the model on the top. The circled graphs form the so-called diamond structure.

Therefore, the family of PD-CGMs, under the model inclusion order, forms a complete, non-distributive lattice.

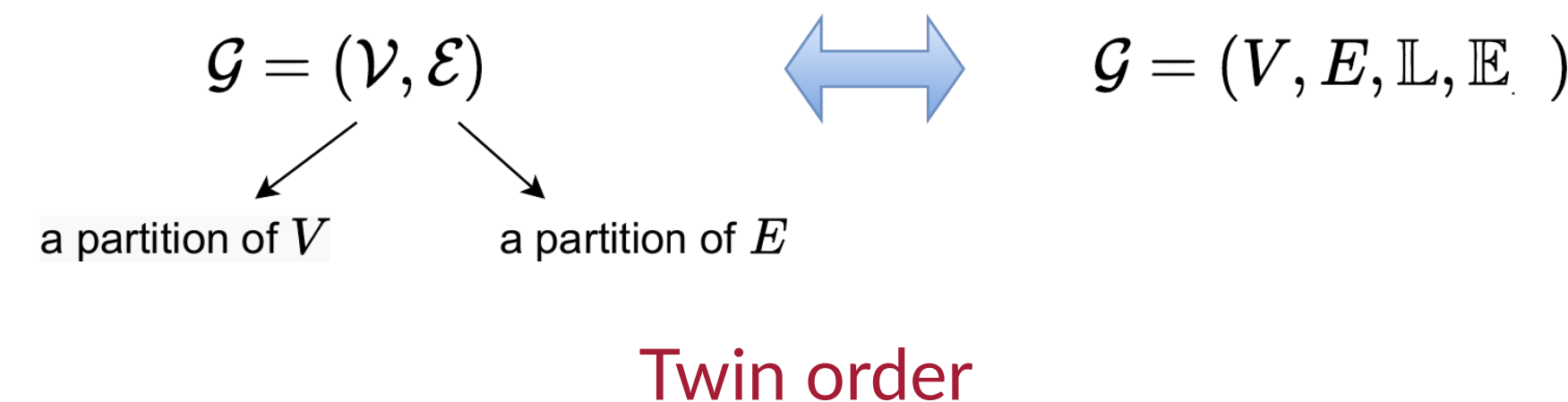
Novel partial order for PD-CGs

The **twin correspondence** $\tau(\cdot)$ is a function of $i \in V$ that is $i + q$ if $i \in L$, and $i - q$ if $i \in R$. Moreover, for $i, j \in V$, $\tau((i, j)) = (\tau(i), \tau(j))$.

We say i, j are **twin vertices** i, j if $\tau(i) = j$ or $i = \tau(j)$, and $(i, j), (k, l)$ are **twin edges** if $\tau(i, j) = (k, l)$ or $(i, j) = \tau(k, l)$.



An alternative and equivalent representation of PD-CGs.



Twin order

For two PD-CGs \mathcal{G} and \mathcal{H} , we say $\mathcal{G} \preceq_\tau \mathcal{H}$ if and only if

- $E_G \subseteq E_H$,
- $\mathbb{L}_G \subseteq \mathbb{L}_H$,
- $\mathbb{E}_G \subseteq \mathbb{E}_H$.

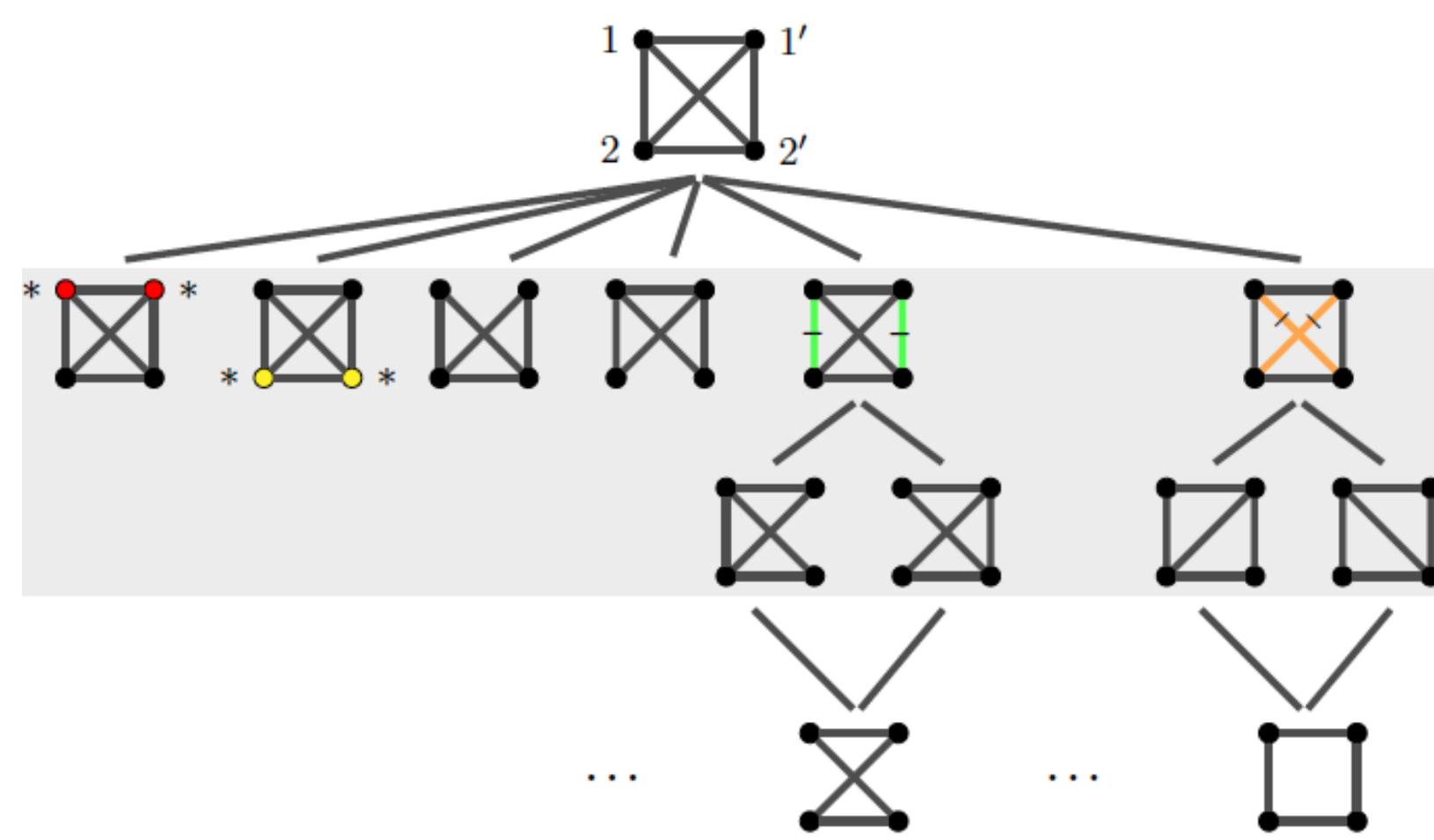
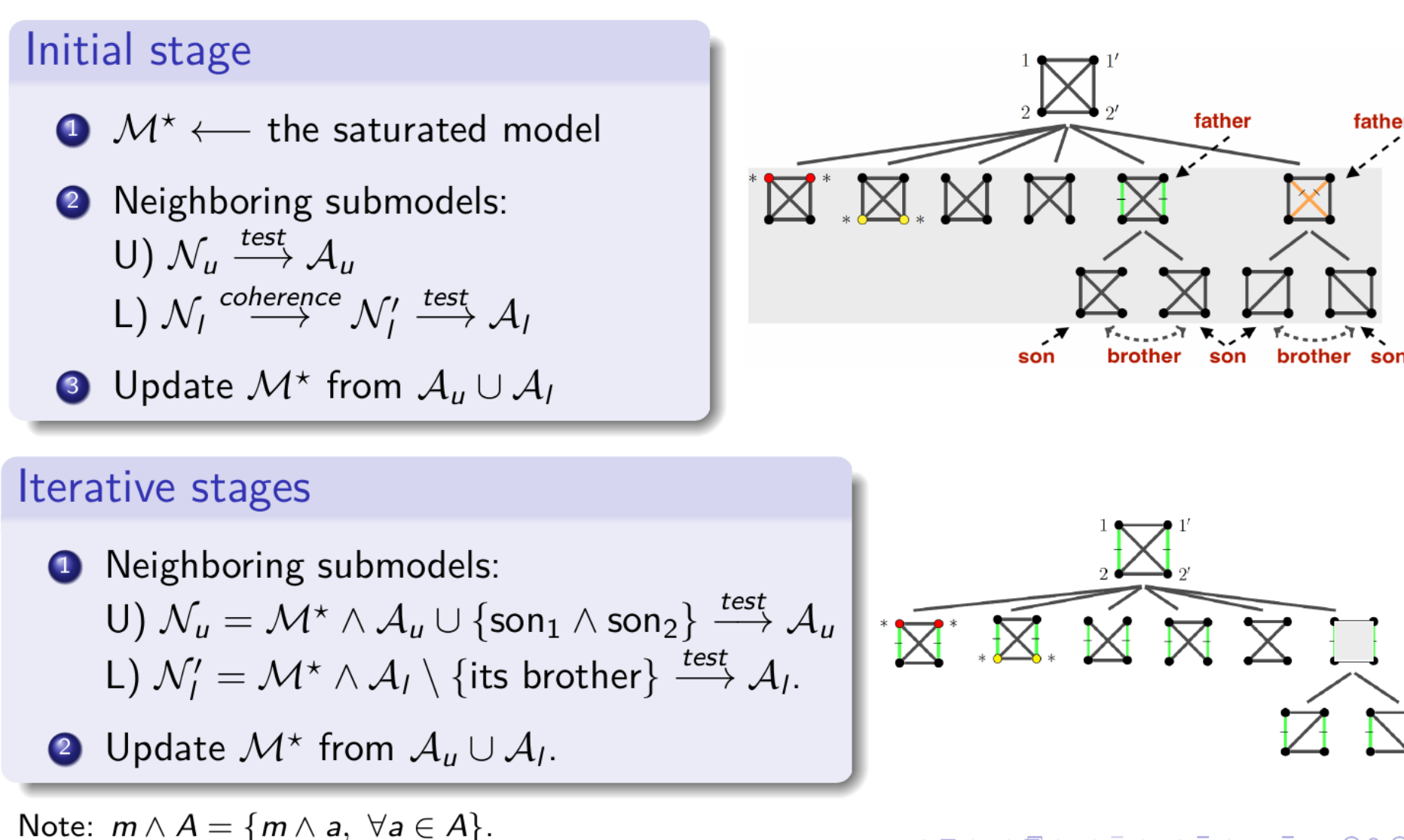


Figure 3. A part of Hasse diagram of the lattice structure of PD-CGs with 4 vertices based on the twin order. The highlighted graphs are the neighbors of the model on the top.

Theorem. The family of PD-CGs under the twin order forms a complete and distributive lattice.

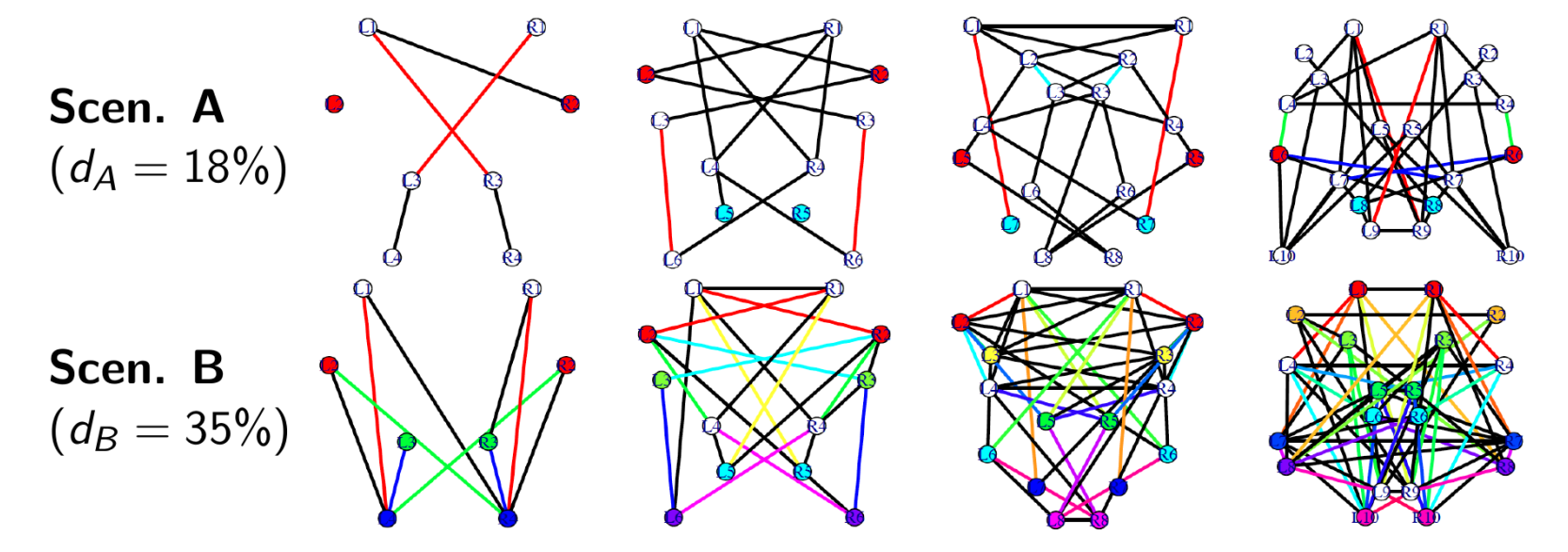
Proposition. For two PD-CGs \mathcal{G}, \mathcal{H} , if $\mathcal{G} \preceq_s \mathcal{H}$ then $\mathcal{G} \preceq_\tau \mathcal{H}$.

Backward elimination stepwise procedure with coherent steps



Numerical experiment

- We generate 100 independent samples with different numbers of variables p varying in $\{8, 12, 16, 20\}$. The recorded results are taken on average over 20 simulated data sets.



Recorded results

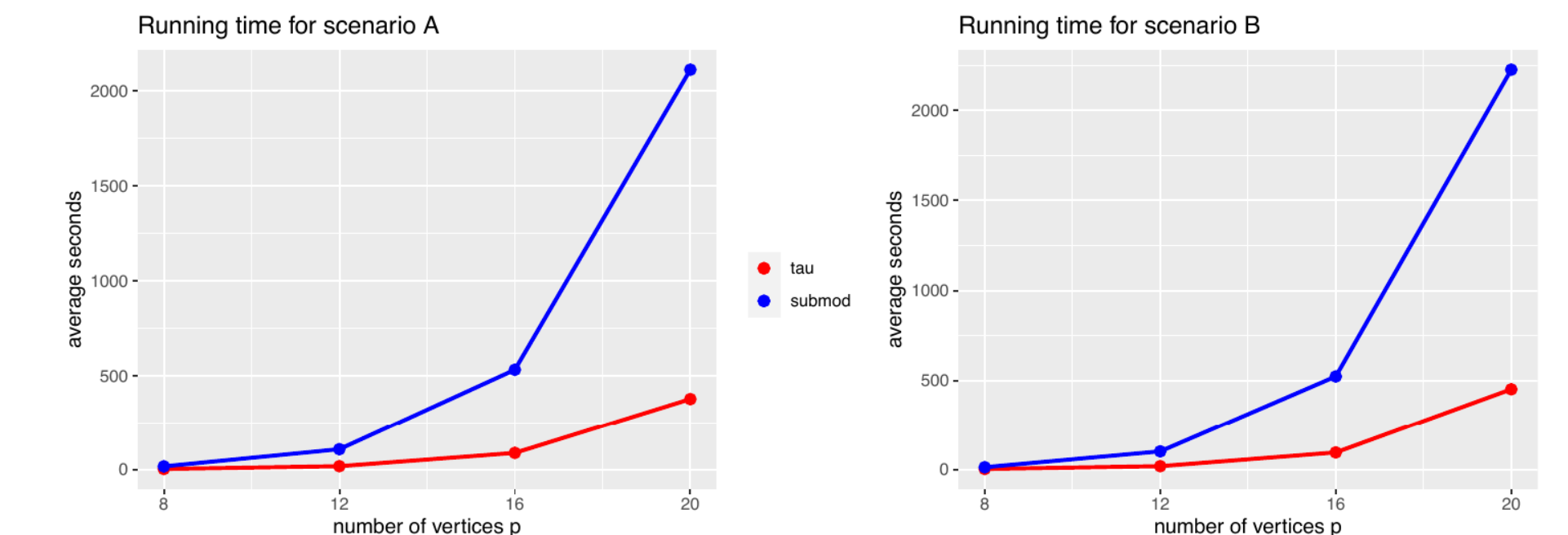


Figure 4. Elapsed time from the backward elimination procedures based on the twin order \preceq_τ (illustrated in red) and the model inclusion \preceq_s (illustrated in blue) of two scenarios A (on the left) and B (on the right).

Table 1. Performance measures of the model selection procedure for the lattice structure equipped by the partial orders \preceq_τ and \preceq_s .

S	p	Order	#edges	Graph structure ePPV _{0%}	eTPR _{0%}	eTNR _{0%}	#sym	Symmetries ePPV _{0%}	eTPR _{0%}	eTNR _{0%}	Time(s)	#models
A	8	\preceq_τ	7(2)	76.68	100.00	91.52	2(1)	41.67	95.00	89.44	4	273
	8	\preceq_s	7(2)	75.41	100.00	91.30	2(1)	46.67	95.00	85.56	17	580
	12	\preceq_τ	17(3)	71.22	97.92	90.37	6(1)	15.99	90.00	87.61	19	1300
	12	\preceq_s	17(3)	70.23	98.75	90.00	5(1)	17.34	90.00	83.91	109	2985
	16	\preceq_τ	27(4)	74.83	88.64	92.70	9(1)	18.53	85.00	89.43	89	4245
B	8	\preceq_τ	28(4)	70.98	87.05	91.48	8(1)	19.32	77.50	84.77	532	10554
	8	\preceq_s	44(8)	64.24	82.21	89.49	16(3)	13.47	70.00	86.18	379	10212
	8	\preceq_τ	46(7)	60.11	78.97	88.04	13(3)	11.97	51.67	80.00	2102	27356
	12	\preceq_τ	11(2)	84.54	89.50	89.72	5(1)	64.08	93.33	92.50	3	264
	12	\preceq_s	11(2)	83.59	89.00	89.44	4(1)	64.83	85.00	85.83	15	486

Concluding remarks:

- The model selection procedure on the twin lattice \preceq_τ is considerably faster than the similar approach on \preceq_s , as shown in Figure 4 and Table 1.
- With $p = 36$, the procedure with the twin order ≈ 7 hours whereas the existing procedure is infeasible.
- The procedure with the twin order tends to perform better when many symmetries are present.

Application on fMRI data

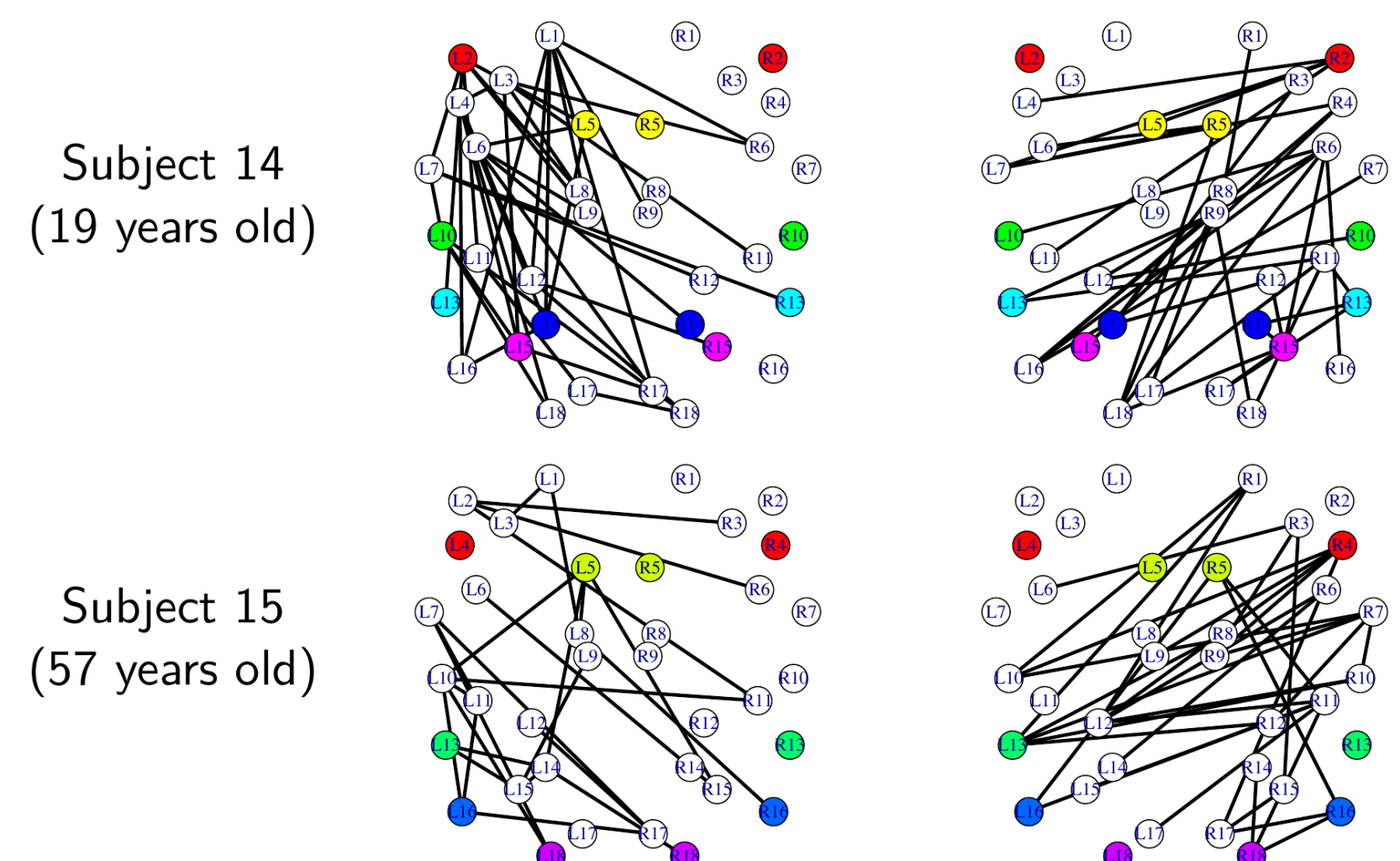


Figure 5. Coloured graphical representations for 36 brain regions in anterior temporal and frontal lobes between two hemispheres.

References

- [1] Davey, B. A. and Priestley, H. A. (2002) *Introduction to lattices and order*. Cambridge University Press.
- [2] Gehrmann, H. (2011) Lattices of graphical Gaussian models with symmetries. *Symmetry*, **3**(3), 653 – 679.
- [3] Hojsgaard, S. and Lauritzen, S. L. (2008) Graphical Gaussian models with edge and vertex symmetries. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **70**(5), 1005 – 1027.
- [4] Ranciat, S., Roverato, A. and Luati, A. (2021) Fused graphical lasso for brain networks with symmetries. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **70**(5), 1299 – 1322.
- [5] Roverato, A. and Nguyen, D. N. (2022) Model inclusion lattice of coloured Gaussian graphical models for paired data. *Proceedings of The 11th International Conference on Probabilistic Graphical Models*, PMLR **186**, 133 – 144.
- [6] Roverato, A. and Nguyen, D. N. Stepwise model search for multiple Gaussian graphical models for paired data (*working paper*).

Contact information

- Dung Ngoc NGUYEN, Postdoctoral Research Fellow.
- Department of Statistical Sciences, University of Padova.
- ngocdung.nguyen@unipd.it
- <https://ngocdung-nguyen.github.io/>
- <https://github.com/NgocDung-NGUYEN>

