

Adolescent paranoia: Prevalence, structure, and causal mechanisms

Jessica C. Bird^{a, d*}, Robin Evans^b, Felicity Waite^{a, d}, Bao S. Loe^c, & Daniel Freeman^{a, d}

Online Supplement 2

This statistical supplement provides further details of the DAGs analysis and the pairwise plots of the causal effects between all variables. The R code for our analysis is available on request.

Directed Acyclic Graphs (DAGs) analysis

The DAGs analysis is described in the main paper. For readers with a statistical background, the analysis used the partition MCMC algorithm of Kuipers and Moffa¹ applied to the covariance matrix of the transformed data. Partition MCMC is a search and score method that samples DAG structures using the Bayesian Gaussian equivalent (BGe) score. BGe scores represent the marginal likelihood of each DAG model under a Wishart prior on the parameters. This formulation implies that the posterior distribution is also Wishart, and the marginal likelihood of the model can be calculated in closed form. This allows calculation, up to a normalising constant, of the posterior probability of each DAG structure. Once the DAGs have been sampled, we are able to obtain posterior samples of the covariance matrix using the posterior inverse Wishart distribution.

Although the analysis followed the approach of Moffa et al.², there were some differences in our approach. Principally, we did not dichotomise the data, but rather mapped the values of each variable to the quantiles of a standard normal distribution. The causal effects we estimate are therefore on the scale of z -scores representing standard deviations in the standardised space. This approach is in line with the *nonparanormal* approach of Liu et al.³ Note that since our data are counts including ties we can only approximate to the multivariate Gaussian distribution. However, this approximation to a continuous distribution is reasonable given that most observations took distinct values.

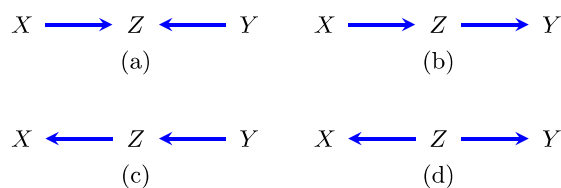
A note on equivalence classes in DAGs

Statistically, a DAG model assumes that the joint distribution of all random variables factorises according to the structure of the graph:

$$P(X_V = x_V) = \prod_{i=1}^k P(X_i = x_i | X_{pa(i)} = x_{pa(i)})$$

In words, the distribution of each variable only depends upon its parents in the graph. Since this factorisation places restrictions on the joint probability distribution of the variables, it is possible to learn the underlying DAG from observed patterns in the data. However, in some cases different causal models imply the same conditional independence relationships between the variables. The simplest case is that the two graphs $X \rightarrow Y$ and $X \leftarrow Y$ cannot be distinguished. More generally, it will occur when two graphs have the same adjacent pairs of nodes, and the same *v-structures*: these are pairs of edges of the form $X \rightarrow Z \leftarrow Y$, where X and Y are not joined by an edge directly. Such DAGs form a *Markov equivalence class* and cannot be distinguished from the patterns of dependence and independence within the data. As a result, it is not

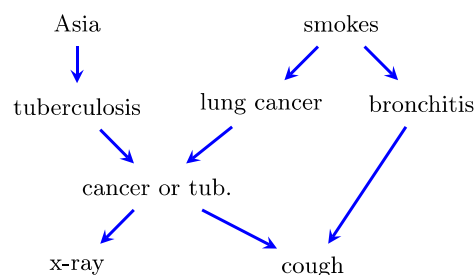
possible to learn the direction of all causal links within a set of variables. As an example, the three graphs (b), (c) and (d) in the figure on the right are all Markov equivalent, while graph (a) is not Markov equivalent to the others because it contains a v-structure⁴.



A note on independence in DAGs

To describe the independences implied by a directed acyclic graph model we need some additional terminology. Given two nodes *X* and *Y*, we say that *Y* is a *descendant* of *X* if there is a sequence of edges directed from *X* to *Y*. For example, $X \rightarrow A \rightarrow Y$ or $X \rightarrow Y$. Otherwise *Y* is a *non-descendant* of *X*. We can think of a descendant as being something that happens ‘later’, in a causal sense, because it is affected by what happens now. Formally, a DAG model implies that each node is independent of its non-descendants once its parents have been controlled for. In other words, once the immediate parents are known, no further information about a node can be gained from its ‘past’.

As an example, consider the DAG on the right, adapted from Lauritzen and Spiegelhalter.⁵ This graph represents the relationship between three diseases (tuberculosis, lung cancer and bronchitis), two risk factors (smoking and a recent visit to Asia), and two symptoms (a shadow on an x-ray and a cough). Following the rule above we can deduce that, under the associated model, the risk of lung cancer is independent of the risk of tuberculosis and bronchitis, once smoking status has been taken into account.



Interpretation of credible intervals

The 90% credible intervals used for causal effects in Table 2 of the main paper only include the sampled graphs in which the relevant causal pathway was present. Hence, they should be interpreted conditionally on the existence of such a pathway. Consider the causal effect of body image on paranoia. A causal pathway was present from body image to paranoia in 46% of graphs, a pathway in the other direction in 45% of graphs, and no pathway in the remaining 9%. Within the 46% of graphs in which the pathway was present, the average causal effect was $z_t = -0.21$, with 90% credible interval from -0.49 to -0.01 . Note that an *overall* credible interval for this causal effect across all graphs would certainly include zero, since this happens in 54% of sampled graphs. Such an interval is potentially misleading, since its values would be unstable if the proportion of graphs in which the pathway were present was close to 95%. Of those graphs containing some causal pathway from body image to paranoia, 44% included the direct edge from body image to paranoia; the strength of the edge in those graphs (now only 20% of total sampled graphs) was $z_d = -0.05$, with 90% credible interval $(-0.15, 0.00)$.

Pairwise causal plots

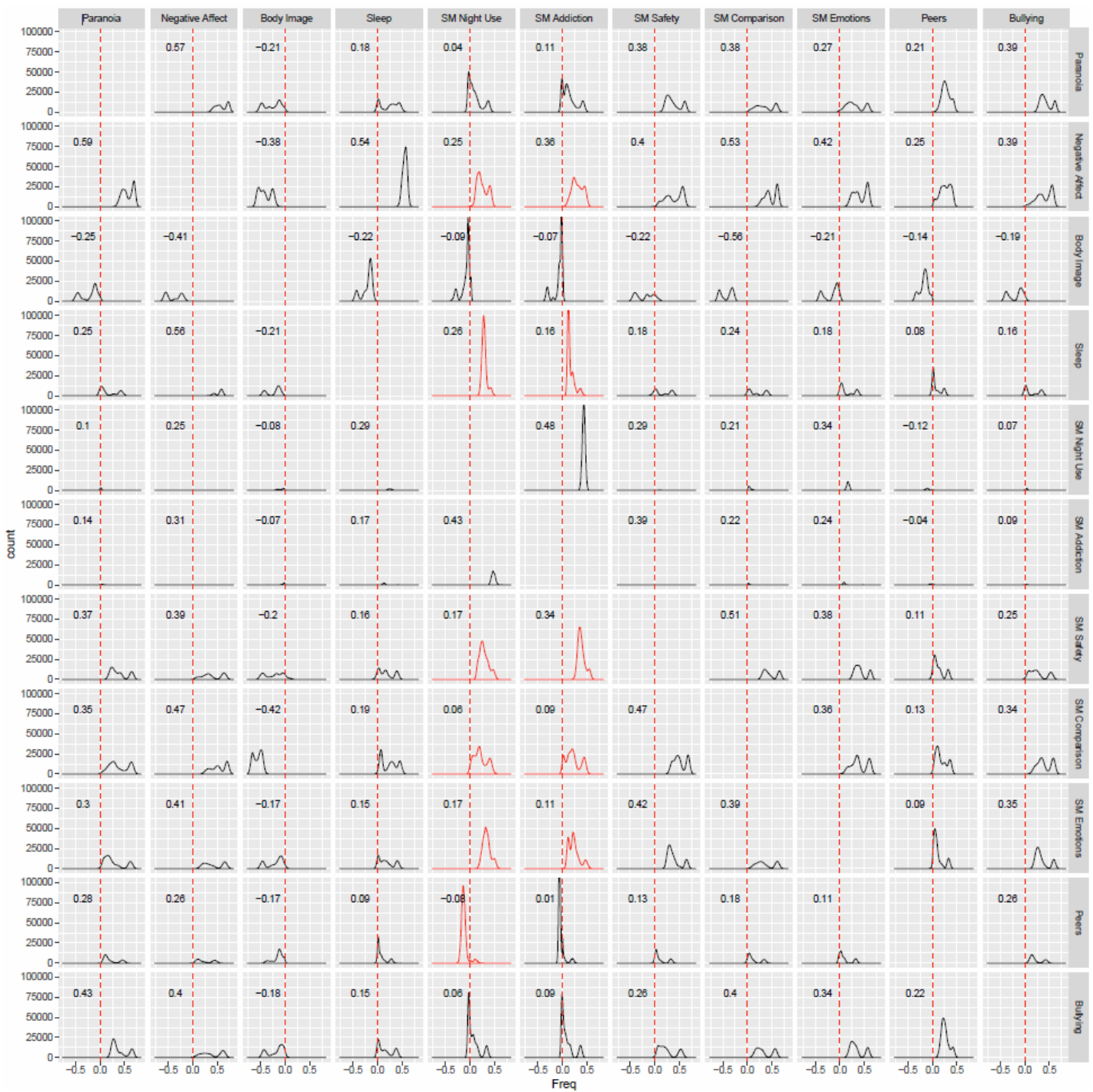


Fig 1. Plots of the average total causal effects for each variable on every other variable. $\bar{\alpha}$ scores of the causal effect is shown on each plot. Red dotted line indicates zero causal effect. Red plots indicate a significant directed causal effect. SM = social media. SM Safety = social media safety-seeking behaviours.

References

- 1 Kuipers J, Moffa G. Partition MCMC for inference on acyclic digraphs. *J Am Stat Assoc* 2017; 112: 282–99.
- 2 Moffa G, Catone G, Kuipers J, et al. Using directed acyclic graphs in epidemiological research in psychosis: An analysis of the role of bullying in psychosis. *Schizophr Bull* 2017; 43: 1273–9.
- 3 Liu H, Lafferty J, Wasserman L. The nonparanormal: Semiparametric estimation of high dimensional undirected graphs. *J Mach Learn Res* 2009; 10: 2295–328.
- 4 Pearl J. *Causality: Models, Reasoning and Inference*. Cambridge, UK: Cambridge University Press, 2009.
- 5 Lauritzen SL, Spiegelhalter DJ. Local computations with probabilities on graphical structures and their application to expert systems. *J R Stat Soc Ser B* 1988; 50: 157–224.