

Supplementary Material for “Graph-structured variable selection with Gaussian Markov random field horseshoe prior”

Marie Denis¹ and Mahlet G. Tadesse²

¹*CIRAD, UMR AGAP Institut, Montpellier, France*

²*Department of Mathematics and Statistics, Georgetown University,
Washington, DC, USA*

This supplement presents additional results for the simulation studies. We report the CPs and widths of the 95% highest posterior density (HPD) intervals when the regression coefficients have varying signs under the non-sequential dependence structure. We also report all the results when the regression coefficients have the same sign under sequential dependence with $p = 1000$ (the results for $p = 140$ are presented in the main paper), and under non-sequential dependence with both $p = 140$ and $p = 1000$. The results using an estimated graph are also presented for the non-sequential simulated setting when the non-zero regression coefficients in a group have the same sign. Finally, in order to show that the choice of the representative vertex from each disjoint subgraph can be arbitrary, we provide the results of repeating the analysis with 20 randomly selected vertices in the sequential dependence scenario.

S1. Results when regression coefficients have varying signs under non-sequential dependence

Table 1 presents the coverage probabilities (CPs) and widths of the 95% HPD intervals for the zero and non-zero regression coefficients under the two

covariance structures considered. For the zero coefficients, the four methods provide similar CPs with all of the 95% HDP intervals containing zero. SS-Ising provides narrower intervals, which, in part, can be explained by the fact that the spike-and-slab prior sets exactly to zero the regression coefficients for non-selected covariates. For the non-zero coefficients, SS-Ising still yields narrower HPD intervals but misses the true parameter values as evidenced by the low CPs. On the other hand, the HS-based approaches yield good CPs with HS-GMRF having relatively narrower intervals. This indicates better accuracy in the estimations obtained with HS-GMRF.

To gain insights into these results, we examine the estimation of all the regression coefficients in one simulated replication (Figure 1). We note that all the methods capture the zero coefficients. For the non-zero coefficients, the HS and SS-Ising approaches tend to select a few representatives among groups of positively correlated or negatively correlated covariates and shrink the coefficients of the non-selected covariates to 0. HS also gives wide HPD intervals, hence leading to good CPs but poor selection performance. On the other hand, the HS-GMRF-based approaches tend to give more accurate coefficient estimation with similar estimates for highly correlated covariates. This explains why the large effect size $\beta_{gk} = 5$ is underestimated. We note some differences when the correlation sign is incorporated or not. For HS-GMRF-nosign, the HPD intervals are relatively wider (as observed in Table 1) and this is more pronounced for coefficients that have opposite signs. Instead, HS-GMRF yields narrower HPD intervals with good coverage and provides fairly accurate estimates for regression coefficients with opposite signs.

For further exploration, we focus on the posterior densities of three non-zero regression coefficients in one of the groups of correlated covariates (Figure 2). For the large effect size, $\beta_{gk} = 5$, HS leads to a bimodal posterior density spanning a wide range of values, resulting in a 95% HPD interval that covers the true value but is too wide and not precise. SS-Ising leads to a multimodal posterior density that also spans a wide range. The HS-GMRF approaches, on the other hand, lead to posterior densities that are narrower and away from 0, but do not capture the true value. For moderate values of $\beta_{gk} = \pm 5/\sqrt{10}$, HS is concentrated around zero and does not capture the true values. SS-Ising displays a similar behavior with a stronger peak around 0. For $\beta_{gk} = 5/\sqrt{10}$, both HS-GMRF approaches lead to narrow posterior densities that are concentrated around the true value, with HS-GMRF-nosign covering a slightly wider range and being less accurate. For $\beta_{gk} = -5/\sqrt{10}$, HS-GMRF-nosign leads to a relatively spread out posterior density peaked

around the average of the β_{gk} values, while HS-GMRF results in a posterior density concentrated around the true value. This indicates that integrating the correlation sign provides a more accurate estimation when the regression coefficients of connected covariates have varying signs.

Table 1: Coverage probability (CP) and width of 95% HPD intervals averaged over the 50 simulated replications when regression coefficients have varying signs under non-sequential dependence.

		CP of 95% HPD			
		Overall	β_j 's = 0	β_j 's \neq 0	
$\Sigma_{g,\text{half}}$	$\rho = 0.5$	HS-GMRF	0.923 (± 0.026)	1.000 (± 0.000)	0.786 (± 0.073)
		HS-GMRF nosign	0.931 (± 0.027)	0.998 (± 0.007)	0.810 (± 0.072)
		HS	0.894 (± 0.037)	1.000 (± 0.000)	0.704 (± 0.102)
		SS-Ising	0.751 (± 0.026)	0.999 (± 0.005)	0.306 (± 0.070)
	$\rho = 0.9$	HS-GMRF	0.928 (± 0.019)	1.000 (± 0.000)	0.800 (± 0.054)
		HS-GMRF nosign	0.922 (± 0.031)	1.000 (± 0.000)	0.782 (± 0.086)
		HS	0.908 (± 0.05)	1.000 (± 0.000)	0.743 (± 0.14)
		SS-Ising	0.773 (± 0.029)	1.000 (± 0.000)	0.365 (± 0.079)
$\Sigma_{g,\text{all}}$	$\rho = 0.5$	HS-GMRF	0.894 (± 0.048)	1.000 (± 0.000)	0.704 (± 0.136)
		HS-GMRF nosign	0.938 (± 0.018)	0.999 (± 0.005)	0.829 (± 0.049)
		HS	0.909 (± 0.041)	1.000 (± 0.000)	0.745 (± 0.116)
		SS-Ising	0.777 (± 0.029)	1.000 (± 0.000)	0.365 (± 0.079)
	$\rho = 0.9$	HS-GMRF	0.963 (± 0.003)	1.000 (± 0.000)	0.898 (± 0.008)
		HS-GMRF nosign	0.945 (± 0.014)	1.000 (± 0.000)	0.847 (± 0.038)
		HS	0.898 (± 0.069)	1.000 (± 0.000)	0.716 (± 0.193)
		SS-Ising	0.813 (± 0.031)	1.000 (± 0.000)	0.475 (± 0.088)

		Width of 95% HPD			
		Overall	β_j 's = 0	β_j 's \neq 0	
$\Sigma_{g,\text{half}}$	$\rho = 0.5$	HS-GMRF	2.047 (± 0.188)	1.900 (± 0.182)	2.311 (± 0.216)
		HS-GMRF nosign	2.712 (± 0.231)	2.446 (± 0.216)	3.193 (± 0.297)
		HS	2.871 (± 0.278)	2.476 (± 0.214)	3.584 (± 0.415)
		SS-Ising	0.656 (± 0.117)	0.249 (± 0.081)	1.389 (± 0.251)
	$\rho = 0.9$	HS-GMRF	2.415 (± 0.248)	2.164 (± 0.228)	2.866 (± 0.31)
		HS-GMRF nosign	3.212 (± 0.284)	2.636 (± 0.237)	4.249 (± 0.432)
		HS	3.255 (± 0.419)	2.483 (± 0.294)	4.646 (± 0.699)
		SS-Ising	0.927 (± 0.181)	0.231 (± 0.132)	2.181 (± 0.438)
$\Sigma_{g,\text{all}}$	$\rho = 0.5$	HS-GMRF	1.306 (± 0.129)	1.275 (± 0.136)	1.362 (± 0.127)
		HS-GMRF nosign	2.439 (± 0.216)	2.208 (± 0.199)	2.854 (± 0.266)
		HS	3.031 (± 0.304)	2.559 (± 0.222)	3.880 (± 0.477)
		SS-Ising	0.772 (± 0.159)	0.233 (± 0.101)	1.741 (± 0.363)
	$\rho = 0.9$	HS-GMRF	1.725 (± 0.241)	1.580 (± 0.214)	1.987 (± 0.313)
		HS-GMRF nosign	3.552 (± 0.366)	2.936 (± 0.353)	4.661 (± 0.462)
		HS	3.651 (± 0.750)	2.737 (± 0.564)	5.297 (± 1.148)
		SS-Ising	1.153 (± 0.218)	0.169 (± 0.108)	2.923 (± 0.547)

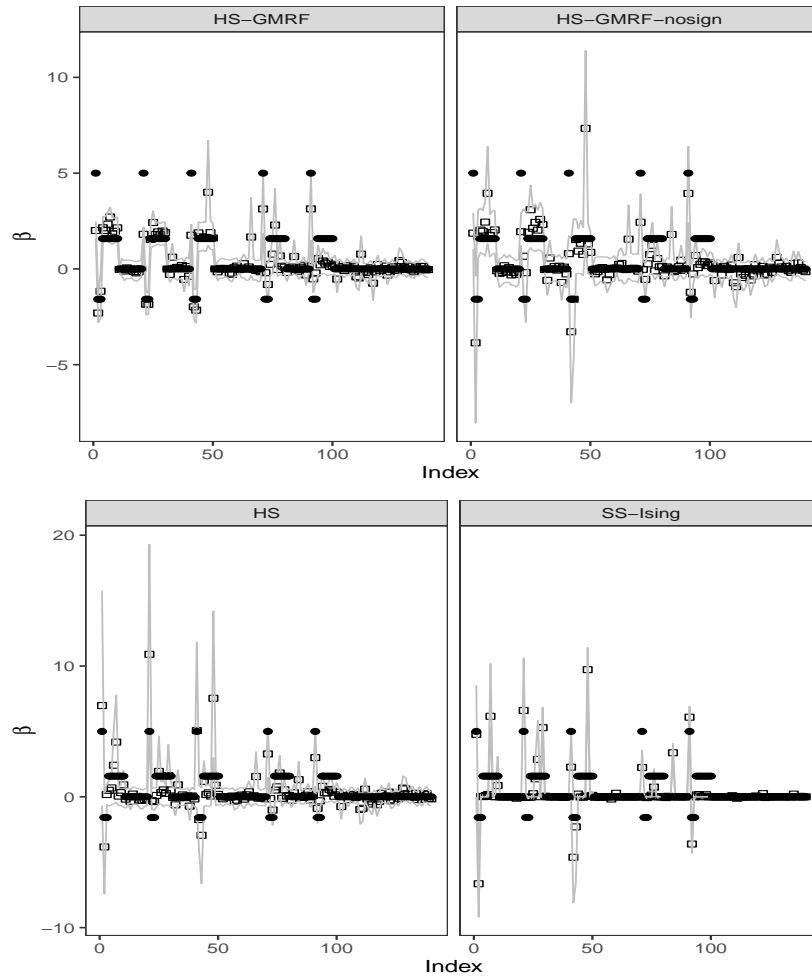


Figure 1: True coefficients (filled black circles) and estimated coefficients (unfilled squares) along with 95% HPD intervals (solid gray lines) for all covariates in one simulated replication with covariance structure $\Sigma_{g,\text{half}}$ and $\rho = 0.9$.

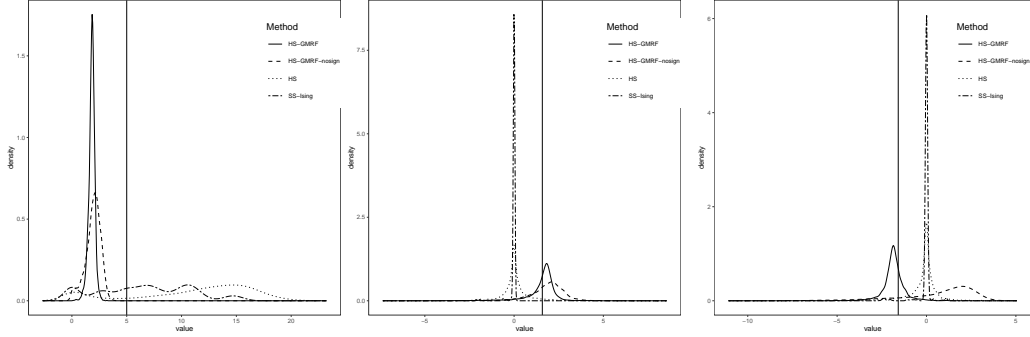


Figure 2: Posterior densities for three non-zero regression coefficients with covariance structure $\Sigma_{g,\text{half}}$ and $\rho = 0.9$. The black vertical bar corresponds to the true β_j value.

S2. Results for sequential dependence

The results with $p = 1000$ covariates are the same as those with $p = 140$ covariates. HS-GMRF has significantly larger MCC, and smaller MSE and MSPE (Table 2). On average, it has smaller coverage and narrower HPD interval width for true signals compared to HS (Table 3).

Table 2: Average MCC, MSE and MSPE (with SE) over 50 simulated replications under sequential dependence with $p = 1000$ covariates.

	MCC	MSE	MSPE
HS-GMRF	0.695 (± 0.050)	0.051 (± 0.006)	128.194 (± 31.521)
HS	0.080 (± 0.020)	0.325 (± 0.140)	200.231 (± 56.470)
SS-Ising	0.206 (± 0.047)	0.204 (± 0.000)	287.419 (± 68.739)

Table 3: Coverage probability (CP) and width of 95% HPD intervals averaged over the 50 simulated replications under sequential dependence with $p = 1000$ covariates.

	CP of 95% HPD		
	Overall	β_j 's = 0	β_j 's \neq 0
HS-GMRF	0.915 (\pm 0.020)	0.997 (\pm 0.007)	0.588 (\pm 0.088)
HS	0.980 (\pm 0.012)	1.000 (\pm 0.000)	0.901 (\pm 0.060)
SS-Ising	0.837 (\pm 0.009)	0.996 (\pm 0.002)	0.197 (\pm 0.041)

	Width of 95% HPD		
	Overall	β_j 's = 0	β_j 's \neq 0
HS-GMRF	0.408 (\pm 0.116)	0.272 (\pm 0.112)	0.952 (\pm 0.163)
HS	1.947 (\pm 0.172)	1.473 (\pm 0.120)	3.841 (\pm 0.441)
SS-Ising	0.445 (\pm 0.090)	0.271 (\pm 0.076)	1.144 (\pm 0.216)

S3. Results when regression coefficients have the same sign under non-sequential dependence

S3.1. Results for $p = 140$ covariates

The HS-GMRF-based methods have much better performance than HS and SS-Ising, in terms of selection, estimation and prediction for all values of ρ and covariance structures considered (Table 4). The MCC values with the HS-GMRF approaches are higher and become closer to 1 under $\Sigma_{g,\text{all}}$ when more covariates are connected. The MSEs for the regression coefficient estimates are both substantially smaller and less variable with the HS-GMRF approaches. In terms of prediction, the HS-GMRF-based approaches provide smaller MSPEs. The HS and SS-Ising methods lead to higher MSPEs. As observed with the simulations under regression coefficient setting 1, where effect sizes have varying signs in the same group, HS has the worst performance in terms of MCC, but has smaller MSEs and MSPEs compared to SS-Ising.

Table 5 provides the coverage probabilities (CP) and widths of the 95% HPD intervals. The HS-GMRF-based approaches yield both higher CPs and narrower widths compared to HS and SS-Ising. As observed with simulation setting 1 for the regression coefficients, SS-Ising has the lowest CP and the

narrowest width under all covariance structures and correlation levels. As noted in the main paper, this is due to SS-Ising missing most of the non-zero regression coefficients.

Table 4: Average MCC, MSE and MSPE (with SE) over 50 simulated replications when regression coefficients have the same sign for $p = 140$.

		MCC	MSE	MSPE
$\Sigma_{g,\text{half}}$ $\rho = 0.5$	HS-GMRF	0.717 (± 0.018)	0.515 (± 0.055)	96.198 (± 13.892)
	HS-GMRF-nosign	0.717 (± 0.018)	0.517 (± 0.055)	96.363 (± 13.695)
	HS	0.247 (± 0.036)	1.049 (± 0.215)	119.421 (± 16.388)
	SS-Ising	0.322 (± 0.044)	1.383 (± 0.252)	139.804 (± 25.492)
$\Sigma_{g,\text{half}}$ $\rho = 0.9$	HS-GMRF	0.670 (± 0.047)	0.550 (± 0.125)	88.651 (± 13.844)
	HS-GMRF-nosign	0.669 (± 0.045)	0.549 (± 0.125)	88.567 (± 13.743)
	HS	0.210 (± 0.044)	2.056 (± 0.627)	100.294 (± 16.893)
	SS-Ising	0.321 (± 0.047)	2.259 (± 0.517)	110.107 (± 18.688)
$\Sigma_{g,\text{all}}$ $\rho = 0.5$	HS-GMRF	0.983 (± 0.016)	0.422 (± 0.021)	69.03 (± 10.441)
	HS-GMRF-nosign	0.986 (± 0.014)	0.422 (± 0.022)	69.077 (± 10.526)
	HS	0.260 (± 0.022)	1.236 (± 0.338)	118.334 (± 22.655)
	SS-Ising	0.351 (± 0.057)	1.660 (± 0.344)	129.761 (± 23.275)
$\Sigma_{g,\text{all}}$ $\rho = 0.9$	HS-GMRF	0.967 (± 0.036)	0.430 (± 0.046)	56.129 (± 9.099)
	HS-GMRF-nosign	0.970 (± 0.023)	0.431 (± 0.047)	56.150 (± 9.052)
	HS	0.198 (± 0.049)	3.116 (± 0.973)	76.679 (± 13.477)
	SS-Ising	0.350 (± 0.049)	3.016 (± 0.664)	74.714 (± 13.19)

Table 5: Coverage probability (CP) and width of 95% HPD intervals (with SE) averaged over the 50 simulated replications when regression coefficients have the same sign for $p = 140$.

			CPs of 95% HPD	Width of 95% HPD
$\Sigma_{g,\text{half}}$	$\rho = 0.5$	HS-GMRF	0.915 (± 0.025)	2.004 (± 0.188)
		HS-GMRF-nosign	0.916 (± 0.024)	2.004 (± 0.195)
		HS	0.883 (± 0.049)	2.784 (± 0.300)
		SS-Ising	0.756 (± 0.033)	0.656 (± 0.171)
	$\rho = 0.9$	HS-GMRF	0.929 (± 0.019)	2.422 (± 0.257)
		HS-GMRF-nosign	0.930 (± 0.021)	2.428 (± 0.258)
		HS	0.915 (± 0.038)	3.356 (± 0.366)
		SS-Ising	0.784 (± 0.028)	0.983 (± 0.202)
$\Sigma_{g,\text{all}}$	$\rho = 0.5$	HS-GMRF	0.900 (± 0.035)	1.344 (± 0.114)
		HS-GMRF-nosign	0.900 (± 0.036)	1.333 (± 0.116)
		HS	0.890 (± 0.046)	2.906 (± 0.285)
		SS-Ising	0.766 (± 0.031)	0.729 (± 0.127)
	$\rho = 0.9$	HS-GMRF	0.964 (± 0.002)	1.691 (± 0.201)
		HS-GMRF-nosign	0.964 (± 0.001)	1.708 (± 0.204)
		HS	0.912 (± 0.063)	3.746 (± 0.667)
		SS-Ising	0.812 (± 0.033)	1.154 (± 0.216)

S3.2. Results with $p = 1000$ covariates

In order to investigate the scalability of the approach, we considered simulations with $p = 1000$ covariates and $\Sigma_{g,\text{half}}$. The results were essentially the same as those with $p = 140$ covariates. As shown in Table 6, the proposed HS-GMRF approach leads to significantly larger MCC and significantly smaller MSE for the regression parameters compared to HS or SS-Ising. In terms of MSPE, HS-GMRF has the smallest value when $\rho = 0.5$ and gives a similar value as HS when $\rho = 0.9$. In both cases, SS-Ising has a significantly larger MSPE.

Table 6: Average MCC, MSE and MSPE (with SE) over 50 simulated replications when regression coefficients have the same sign for $p = 1000$.

		MCC	MSE	MSPE
$\rho = 0.5$	HS-GMRF	0.569 (± 0.011)	0.425 (± 0.041)	526.408 (± 80.977)
	HS	0.090 (± 0.005)	2.173 (± 0.287)	667.806 (± 112.633)
	SS-Ising	0.106 (± 0.036)	0.550 (± 0.008)	1295.870 (± 237.924)
$\rho = 0.9$	HS-GMRF	0.552 (± 0.039)	0.451 (± 0.060)	518.441 (± 76.362)
	HS	0.086 (± 0.011)	5.818 (± 0.786)	518.217 (± 71.864)
	SS-Ising	0.139 (± 0.033)	0.536 (± 0.007)	904.302 (± 139.973)

Table 7 shows that HS-GMRF retains the same values for its coverage probability and HPD interval width for $p = 1000$ and $p = 140$. However, HS drops in its coverage probability and its interval width. This can be explained by the increased number of correlated variables and the fact that HS selects one representative from the correlated set and shrinks all the other parameters to zero.

Table 7: Coverage probability (CP) and width of 95% HPD intervals (with SE) averaged over the 50 simulated replications with same sign regression coefficients for $p = 1000$.

		CPs of 95% HPD	Width of 95% HPD
$\rho = 0.5$	HS-GMRF	0.905 (± 0.008)	2.005 (± 0.275)
	HS	0.835 (± 0.016)	1.935 (± 0.453)
	SS-Ising	0.818 (± 0.011)	0.641 (± 0.214)
$\rho = 0.9$	HS-GMRF	0.908 (± 0.011)	2.239 (± 0.528)
	HS	0.826 (± 0.019)	1.594 (± 0.618)
	SS-Ising	0.838 (± 0.008)	1.065 (± 0.213)

S4. Results for estimated non-sequential dependence

We use the graphical lasso method implemented in the R package `CVglasso` to estimate the precision matrix with the penalty parameter selected using

cross-validation (Friedman, Hastie and Tibshirani, 2008). We compare the results using the true simulated graphs and the estimated ones. For moderate correlation the estimated graphs tend to have less edges than the true simulated graphs, while for high correlation the estimated graphs tend to connect all variables within a group of correlated variables resulting in a higher number of edges than the true simulated graphs.

Table 8 reports the mean MCC, MSE and MSPE along with their standard errors over the 50 simulated replications when true and estimated graphs are considered. As HS does not integrate any structure information, the results are the same for the true and estimated graphs. When the correlation ρ is moderate, the structure from the underestimated graph is less informative than the true one and results in slightly lower and more variable MCC values for the HS-GMRF based methods, but does not appear to affect the SS-Ising results. Conversely, under high correlation, the overestimated graph leads to improved variable selection compared to using the true one for the HS-GMRF-based methods, but gives lower MCC values for SS-Ising. In terms of MSE, when ρ is moderate, the results with the true or estimated graphs are comparable for each method. However with high ρ , the estimation with HS-GMRF-nosign is less accurate. With respect to predictive performance, the use of an underestimated graph slightly increases the MSPE for the HS-GMRF based approaches. Using an overestimated graph does not affect HS-GMRF, but the predictive performance of HS-GMRF-nosign and SS-Ising decrease, both leading to higher MSPE than HS, which ignores the covariate structure. This can be explained by the fact that HS-GMRF-nosign suffers from poor estimation when the regression coefficients of connected variables have different signs, thus leading to poor prediction.

Table 9 presents the CPs and widths of the 95% HPD intervals under the true and estimated graphs. For moderate correlation between the predictors, the true and estimated graphs provide similar results for all methods. For high correlation, the use of an overestimated graph impacts the results of the HS-GMRF-based approaches, leading to smaller CPs and narrower intervals. This emphasizes that an overestimation of the number of edges in the graph leads to an oversmoothing of the regression coefficient estimates for connected variables.

Table 8: Average MCC, MSE and MSPE (with SE) over 50 simulated replications using the true and estimated graphs.

		MCC	MSE	MSPE
True graph $\rho = 0.5$	HS-GMRF	0.708 (± 0.018)	0.513 (± 0.067)	94.871 (± 13.632)
	HS-GMRF-nosign	0.624(± 0.034)	0.728(± 0.155)	122.188(± 21.609)
	HS	0.240(± 0.041)	1.009(± 0.200)	126.252(± 19.657)
	SS-Ising	0.323(± 0.054)	1.386(± 0.204)	149.294(± 27.384)
Estimated graph $\rho = 0.5$	HS-GMRF	0.637 (± 0.076)	0.599 (± 0.133)	108.924 (± 24.509)
	HS-GMRF-nosign	0.567(± 0.078)	0.83(± 0.213)	138.111(± 32.206)
	HS	0.234(± 0.040)	1.007(± 0.200)	126.356(± 19.622)
	SS-Ising	0.319(± 0.047)	1.405(± 0.198)	151.844(± 28.532)
True graph $\rho = 0.9$	HS-GMRF	0.668 (± 0.046)	0.541 (± 0.089)	84.954 (± 14.485)
	HS-GMRF-nosign	0.444(± 0.117)	1.038(± 0.259)	99.123(± 17.694)
	HS	0.219(± 0.038)	2.243(± 0.551)	95.219(± 19.279)
	SS-Ising	0.312(± 0.048)	2.359(± 0.437)	109.387(± 23.713)
Estimated graph $\rho = 0.9$	HS-GMRF	0.728 (± 0.037)	0.453 (± 0.045)	82.196 (± 13.672)
	HS-GMRF-nosign	0.723(± 0.027)	1.621(± 0.128)	124.086(± 19.966)
	HS	0.215(± 0.044)	2.243(± 0.570)	95.010(± 18.973)
	SS-Ising	0.252(± 0.078)	3.189(± 0.809)	173.85(± 34.929)

Table 9: Coverage probability (CP) and width of 95% HPD intervals (with SE) averaged over the 50 simulated replications using the true and estimated graphs.

		CP of 95% HPD	width of 95% HPD
True graph $\rho = 0.5$	HS-GMRF	0.923(± 0.026)	2.047(± 0.188)
	HS-GMRF-nosign	0.931(± 0.027)	2.712(± 0.231)
	HS	0.894(± 0.037)	2.871(± 0.278)
	SS-Ising	0.751(± 0.026)	0.656(± 0.117)
Estimated graph $\rho = 0.5$	HS-GMRF	0.899(± 0.036)	2.217(± 0.327)
	HS-GMRF-nosign	0.919(± 0.036)	2.855(± 0.312)
	HS	0.892(± 0.035)	2.877(± 0.280)
	SS-Ising	0.748(± 0.024)	0.655(± 0.126)
True graph $\rho = 0.9$	HS-GMRF	0.928(± 0.019)	2.415(± 0.248)
	HS-GMRF-nosign	0.922(± 0.031)	3.212(± 0.284)
	HS	0.908(± 0.050)	3.255(± 0.419)
	SS-Ising	0.773(± 0.029)	0.927(± 0.181)
Estimated graph $\rho = 0.9$	HS-GMRF	0.757(± 0.050)	1.330(± 0.167)
	HS-GMRF-nosign	0.743(± 0.025)	1.707(± 0.237)
	HS	0.906(± 0.049)	3.256(± 0.419)
	SS-Ising	0.697(± 0.074)	2.158(± 0.442)

S5. Results with randomly selected representative vertex from each disjoint subgraph

In order to show that the choice of the representative vertex from each disjoint subgraph can be arbitrary, we repeated 20 times the analysis for the sequential dependence scenario with $p = 140$, each time picking a different randomly selected vertex from each disjoint subgraph. The estimated regression coefficients for each covariate with the different representative vertices are essentially the same, as shown in Figure 3. The average MCC over the 20 repetitions was 0.754 with a standard error of 0.032, again demonstrating that the choice of the representative vertex can be arbitrary and does not affect the results.

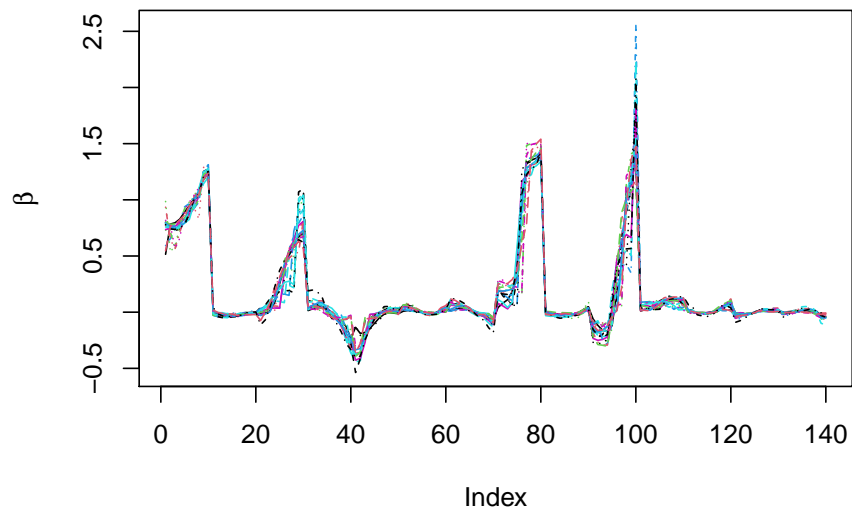


Figure 3: Estimated regression coefficients for $p = 140$ covariates under sequential dependence using different randomly selected representative vertices.

References

Friedman, J., Hastie, T., Tibshirani, R., 2008. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* 9, 432–441.