Supplemental Material.

A comparative study on high-dimensional Bayesian regression with binary predictors.

1. Additional simulations

Due to the comparative nature of the paper, many challenging scenarios may be taken into account in the simulation study. In the paper we focus on those results achieved from the simulations conducted running the set-up used in Griffin and Brown [1]. Following the suggestion of the reviewer, we propose also four very challenging scenarios. Specifically, the supplemental simulations are conducted for n = 100 and they differ for the number of predictors $p = \{200, 5000\}$ and for values of $\beta^* = \{0.1, 0.5\}$. To these scenarios correspond very low values of signal-to-noise ratio. We see that the results confirm the performance of the approaches depicted by those scenarios in which the non-zero coefficients in the regressions are very low and the ratio n/p is also very low.

1.1. Predictive performance

Table 1 shows the performance of the approaches in terms of prediction capacity, i.e. the values of PMSE.

[Table 1 about here.]

1.2. Variable selection performance

We report in Table 2 the median value of TND (in the 50 MCMC runs) for each Bayesian regression, different values of $p = \{100, 5000\}$ and different values of non-zero simulated coefficients $\beta^* = \{0.1, 0.5\}$.

The accuracy in correctly estimating the non-zero regression coefficients are shown in Tables 1 and 4 reporting the FDR and the S measures for the different values of β^* and $p = \{200, 5000\}$.

[Table 2 about here.]

[Table 3 about here.]

[Table 4 about here.]

2. Convergence studies

We provide the Geweke's convergence diagnostic [2] and the Monte Carlo standard errors [3] for the actual non-zero coefficients of the simulation for a subset of selected scenarios characterized by different values of n, p and β^* .

Specifically, for the Geweke's diagnostic we report the test for equality of the means of the first and last part of a Markov chain (the first 10% and the last 20%). If the samples are drawn from the stationary distribution of the chain, the two means are equal and Geweke's statistic (G statistics) has an asymptotically standard normal distribution. The test statistic is a standard Z-score and is calculated under the assumption that the two parts of the chain are asymptotically independent. We use the R function geweke.diag of the package *coda*. Monte Carlo standard errors are calculated for each estimate of the actual non-zero coefficient providing an indication of the variability expected in the estimate. Standard errors of the estimates are calculated using batch means. To compute the Monte Carlo standard errors for our simulations we use the R function mcse.mat of the package *mcmcse*. A graphical representation of the value of the estimates across iterations (trace plot) is also derived.

All results refer to one randomly selected initialization, with total number of MCMC iterations equals to T = 1000 and thinning of the chains t = 100.

2.1. $n=100, p=200, \beta^{\star}=1$

Actual values of the coefficients are: $\beta_1, \beta_{41}, \beta_{81}, \beta_{121}, \beta_{161} = -1$ whereas $\beta_{21}, \beta_{61}, \beta_{101}, \beta_{141}, \beta_{181} = 1$.

Geweke Diagnostic

G statistic represents the Z-score for a test of equality of means between the first and last parts of the chain (H_0) . A G statistic less than $|z_{\alpha/2}|$, where α is the specific significance level for the test, indicates that data support H_0 , so samples are drawn from the stationary distribution of the chain.

G statistics					
	BL	HS	NG	BR	
β_1	1.12	1.24	-0.91	-1.01	
β_{21}	0.48	-3.54	0.02	1.27	
β_{41}	-0.90	-2.31	-0.63	1.15	
β_{61}	-3.15	-0.19	0.73	0.42	
β_{81}	1.69	0.91	-1.02	-0.01	
β_{101}	1.18	0.57	-0.29	-0.20	
β_{121}	1.36	0.18	0.05	0.83	
β_{141}	0.92	-0.73	-0.90	-1.33	
β_{161}	1.66	0.55	0.20	0.54	
β_{181}	-1.27	0.60	-3.21	-1.38	

Fixing a significance level $\alpha = 0.01$ we can see that there are few coefficients not reaching convergence, specifically β_{21} for Horseshoe regression, β_{61} for the Bayesian Lasso and β_{181} for the Normal-Gamma Regression.

Monte Carlo standard errors

Estimates of the actual non-zero coefficients are reported with their associated standard error (in parenthesis).

Estimates (standard error)					
	BL	HS	NG	BR	
β_1	-0.16 (0.01)	-0.10 (0.01)	-0.15 (0.01)	-0.31(0.02)	
β_{21}	0.17(0.01)	0.13(0.01)	0.17(0.01)	$0.40 \ (0.02)$	
β_{41}	-0.06(0.01)	-0.02(0.00)	-0.07(0.01)	-0.12(0.01)	
β_{61}	0.05~(0.01)	$0.02 \ (0.00)$	0.06(0.01)	0.10(0.01)	
β_{81}	-0.28(0.02)	-0.53 (0.03)	-0.32(0.02)	-0.59(0.02)	
β_{101}	0.05~(0.01)	$0.03\ (0.00)$	0.06(0.01)	0.05~(0.01)	
β_{121}	-0.98(0.02)	-1.44(0.03)	-0.98(0.02)	-1.26(0.01)	
β_{141}	0.09(0.01)	$0.07 \ (0.01)$	0.10(0.01)	0.16(0.02)	
β_{161}	-0.67(0.02)	-1.05(0.02)	-0.69(0.02)	-0.80(0.02)	
β_{181}	$0.11 \ (0.01)$	$0.06\ (0.01)$	0.12(0.01)	$0.22 \ (0.02)$	

All the standard errors are small even if many of the estimated coefficients are quite far from their actual values.

Graphical representation of the estimate vs iteration

The trace plots of a set of selected coefficients are reported. Specifically, the estimates of β_1 , β_{81} and β_{161} across iterations are shown for all the regression models. Actual values of the three displayed coefficients are: β_1 , β_{81} , $\beta_{161} = -1$ (depicted with a red line in each graph).



2.2. $n=100, p=200, \beta^{\star}=5$

Actual values of the coefficients are: $\beta_1, \beta_{41}, \beta_{81}, \beta_{121}, \beta_{161} = -5$ whereas $\beta_{21}, \beta_{61}, \beta_{101}, \beta_{141}, \beta_{181} = 5$.

Geweke Diagnostic

G statistic represents the Z-score for a test of equality of means between the first and last parts of the chain (H_0) . A G statistic less than $|z_{\alpha/2}|$, where α is the specific significance level for the test, indicates that data support H_0 , so samples are drawn from the stationary distribution of the chain.

	G statistics					
	BL	HS	NG	BR		
β_1	0.92	0.88	1.78	0.95		
β_{21}	-0.88	-0.89	-0.75	-0.36		
β_{41}	0.90	1.00	0.82	1.01		
β_{61}	-0.64	-0.30	-0.31	0.02		
β_{81}	0.52	0.88	0.90	0.97		
β_{101}	-1.24	-0.58	-0.49	-0.66		
β_{121}	1.26	1.10	0.90	0.45		
β_{141}	-1.04	-1.17	-0.81	-0.38		
β_{161}	0.91	0.41	0.81	0.32		
β_{181}	-0.71	-0.85	-1.05	-0.81		

Fixing a significance level $\alpha = 0.01$, all the coefficients reach the convergence.

Monte Carlo standard errors

Estimates of the actual non-zero coefficients are reported with their associated standard error (in parenthesis).

Estimates (standard error)					
	BL	$_{ m HS}$	NG	BR	
β_1	-5.12(0.04)	-5.37(0.02)	-5.22(0.05)	-4.72(0.06)	
β_{21}	4.40(0.03)	4.62(0.03)	4.54(0.05)	4.10(0.05)	
β_{41}	-4.58(0.02)	-4.61(0.02)	-4.57(0.04)	-4.50(0.03)	
β_{61}	4.61(0.04)	5.08(0.02)	4.89(0.04)	3.93(0.06)	
β_{81}	-4.81(0.04)	-5.28(0.03)	-5.02(0.04)	-4.28(0.07)	
β_{101}	3.74(0.02)	4.19(0.03)	3.93(0.02)	3.40(0.05)	
β_{121}	-5.23(0.03)	-5.40(0.02)	-5.29(0.03)	-4.78(0.05)	
β_{141}	3.55(0.04)	4.18(0.04)	3.86(0.03)	3.12(0.06)	
β_{161}	-4.08(0.03)	-4.61(0.03)	-4.30(0.04)	-3.40(0.06)	
β_{181}	4.96(0.04)	5.24(0.03)	5.09(0.03)	4.52(0.06)	

The standard errors are small for almost all the estimated coefficients. The Bayesian Ridge presents the highest standard errors, and estimated coefficients that are the farest from their actual values with respect to all other Bayesian approaches.

Graphical representation of the estimate vs iteration

The trace plots of a set of selected coefficients are reported. Specifically, the estimates of β_1 , β_{81} and β_{161} across iterations are shown for all the regression models. Actual values of the three displayed coefficients are: β_1 , β_{81} , $\beta_{161} = -5$ (depicted with a red line in each graph).



2.3. $n=100, p=200, \beta^{\star} = \{1, 5, 10\}$

Actual values of the coefficients are: $\beta_1, \beta_{41}, \beta_2, \beta_3, \beta_4 = 1, \beta_5, \beta_6, \beta_7, \beta_8 = 5$ and $\beta_9, \beta_{10} = 10$.

Geweke Diagnostic

G statistic represents the Z-score for a test of equality of means between the first and last parts of the chain (H_0) . A G statistic less than $|z_{\alpha/2}|$, where α is the specific significance level for the test, indicates that data support H_0 , so samples are drawn from the stationary distribution of the chain.

G statistics					
	BL	HS	NG	BR	
β_1	-0.92	1.07	-0.01	1.26	
β_2	1.67	0.05	3.00	0.64	
β_3	-1.84	-1.63	-1.09	0.42	
β_4	-2.29	1.22	-0.44	0.86	
β_5	-0.02	-1.17	-0.86	-0.84	
β_6	-0.49	-1.20	-1.12	-1.33	
β_7	-0.65	-1.15	-1.03	-1.17	
β_8	-0.17	-0.92	-1.95	-1.80	
β_9	-0.69	-0.91	-1.23	-1.24	
β_{10}	-0.72	-0.88	-1.07	-1.40	

Fixing a significance level $\alpha = 0.01$, all the coefficients reach the convergence except for β_2 in the Normal-Gamma regression.

Monte Carlo standard errors

Estimates of the actual non-zero coefficients are reported with their associated standard error (in parenthesis).

Estimates (standard error)					
	BL	HS	NG	BR	
β_1	0.01(0.01)	0.05 (0.01)	0.04(0.01)	$0.07 \ (0.03)$	
β_2	0.19(0.02)	0.12(0.01)	$0.16 \ (0.02)$	0.22(0.03)	
β_3	1.06(0.02)	1.10(0.02)	$1.12 \ (0.02)$	1.02(0.06)	
β_4	0.47(0.04)	$0.42 \ (0.03)$	$0.39\ (0.03)$	$0.50 \ (0.06)$	
β_5	4.25(0.03)	4.68(0.04)	4.45 (0.03)	4.34(0.07)	
β_6	4.45(0.03)	4.68(0.05)	4.56(0.05)	4.44(0.06)	
β_7	5.08(0.04)	5.41(0.07)	5.29(0.03)	5.07(0.08)	
β_8	4.49(0.03)	4.85(0.05)	4.65(0.02)	4.55(0.05)	
β_9	10.33 (0.05)	10.59(0.12)	$10.47 \ (0.03)$	10.14(0.10)	
β_{10}	9.68(0.05)	9.98(0.07)	9.86(0.03)	9.44(0.05)	

The standard errors are small for almost all the estimated coefficients and the Bayesian Ridge presents the highest values of standard errors. The estimates of the coefficients with high actual values (specifically, from β_5 to β_{10}) are very close to their actual values.

Graphical representation of the estimate vs iteration

The trace plots of a set of selected coefficients are reported. Specifically, the estimates of β_1 , β_5 and β_9 across iterations are shown for all the regression models. Actual values of the three displayed coefficients are: $\beta_1 = 1$, $\beta_5 = 5$ and $\beta_9 = 10$ (depicted with a red line in each graph).



2.4. $n=500, p=1000, \beta^{\star}=1$

Actual values of the coefficients are: $\beta_1, \beta_{201}, \beta_{401}, \beta_{601}, \beta_{801} = -1$ whereas $\beta_{101}, \beta_{301}, \beta_{501}, \beta_{701}, \beta_{901} = 1$.

Geweke Diagnostic

G statistic represents the Z-score for a test of equality of means between the first and last parts of the chain (H_0) . A G statistic less than $|z_{\alpha/2}|$, where α is the specific significance level for the test, indicates that data support H_0 , so samples are drawn from the stationary distribution of the chain.

G statistics					
	BL	HS	NG	BR	
β_1	76.25	1.50	102.17	3.56	
β_{101}	-1.09	1.15	-87.76	-1.29	
β_{201}	2.39	96.64	17.96	59.04	
β_{301}	-9.81	-74.28	0.45	-66.39	
β_{401}	71.70	71.31	1.96	1.56	
β_{501}	-58.35	-4.71	-89.51	-56.00	
β_{601}	4.09	1.20	47.04	18.48	
β_{701}	-104.29	-1.70	-97.69	-1.57	
β_{801}	101.90	1.12	3.62	1.83	
β_{901}	-9.16	-16.01	-3.28	-1.17	

Many convergence problems are detected. Fixing a significance level $\alpha = 0.01$, the Bayesian Lasso provides not converging estimates for all the coefficients except for 2 of them (β_{101} and β_{201}); the Horseshoe regression doesn't reach the convergence for half of the coefficients (β_{201} , β_{301} , β_{401} , β_{501} and β_{901}); the Normal-Gamma regression doesn't converge for almost all the coefficients except for 2 of them (β_{301} and β_{401}); the Ridge regression doesn't reach the convergence for half of the coefficients (β_1 , β_{201} , β_{301} , β_{501} and β_{601}).

Monte Carlo standard errors

Estimates of the actual non-zero coefficients are reported with their associated standard error (in parenthesis).

Estimates (standard error)					
	BL	HS	NG	BR	
β_1	-0.78(0.37)	-1.06(0.08)	-0.88(0.33)	-1.08(0.03)	
β_{101}	0.86(0.03)	0.97(0.01)	0.84(0.19)	0.87 (0.05)	
β_{201}	-0.82(0.09)	-0.79(0.17)	-0.57(0.50)	-0.67(0.39)	
β_{301}	0.74(0.06)	0.65(0.33)	0.86(0.01)	0.68(0.15)	
β_{401}	-0.62(0.45)	-0.76(0.24)	-0.91(0.01)	-0.79(0.06)	
β_{501}	0.86(0.12)	0.94(0.11)	0.71(0.49)	0.87(0.19)	
β_{601}	-0.85(0.02)	-0.87(0.05)	-0.74(0.21)	-0.74 (0.10)	
β_{701}	0.83(0.17)	$0.93 \ (0.07)$	0.80(0.24)	0.88(0.07)	
β_{801}	-0.87(0.53)	-1.10(0.04)	-1.10(0.01)	-1.00(0.06)	
β_{901}	$1.04 \ (0.20)$	$1.14 \ (0.19)$	1.19(0.02)	$1.09\ (0.05)$	

Monte Carlo standard errors show high values for many of the coefficients with respect to all the Bayesian regressions. This supports the convergence problems highlighted by the Geweke diagnostic.

Graphical representation of the estimate vs iteration

The trace plots of a set of selected coefficients are reported. Specifically, the estimates of β_1 , β_{401} and β_{801} across iterations are shown for all the regression models. Actual values of the three displayed coefficients are: $\beta_1, \beta_{401}, \beta_{801} = -1$ (depicted with a red line in each graph).



2.5. $n=500, p=1000, \beta^{\star}=5$

Actual values of the coefficients are: $\beta_1, \beta_{201}, \beta_{401}, \beta_{601}, \beta_{801} = -5$ whereas $\beta_{101}, \beta_{301}, \beta_{501}, \beta_{701}, \beta_{901} = 5$.

Geweke Diagnostic

G statistic represents the Z-score for a test of equality of means between the first and last parts of the chain (H_0) . A G statistic less than $|z_{\alpha/2}|$, where α is the specific significance level for the test, indicates that data support H_0 , so samples are drawn from the stationary distribution of the chain.

G statistics					
	BL	HS	NG	BR	
β_1	2.95	469.21	462.67	2.38	
β_{101}	-422.07	-463.25	0.48	-270.74	
β_{201}	281.76	425.55	473.28	193.20	
β_{301}	-7.82	-1.33	-485.66	-312.79	
β_{401}	514.04	521.64	487.06	1.98	
β_{501}	-508.05	-503.24	-267.81	-1.65	
β_{601}	4.94	3.66	399.52	9.05	
β_{701}	-1.36	-3.68	-466.65	-2.46	
β_{801}	2.51	380.80	360.10	317.01	
β_{901}	-1.15	-0.47	-501.59	-1.42	

Many convergence problems are detected and all the approaches fail in reaching the convergence for almost all the coefficients. Fixing a significance level $\alpha = 0.01$, we can see that the Bayesian Ridge reaches the convergence for half of the coefficients (β_1 , β_{401} , β_{501} , β_{701} and β_{901}).

Monte Carlo standard errors

Estimates of the actual non-zero coefficients are reported with their associated standard error (in parenthesis).

Estimates (standard error)					
	BL	HS	NG	BR	
β_1	-4.61(0.39)	-4.37(0.90)	-4.28(1.08)	-4.60(0.45)	
β_{101}	3.55(1.76)	4.02(0.86)	4.88(0.11)	4.05(0.89)	
β_{201}	-3.85(1.98)	-3.98(1.47)	-3.80(1.92)	-4.35(0.94)	
β_{301}	5.13(0.11)	5.05(0.11)	3.93(1.93)	4.30(1.12)	
β_{401}	-3.31(2.59)	-4.44(0.86)	-2.43(3.05)	-4.91(0.28)	
β_{501}	4.41 (0.81)	4.02(1.31)	3.94(1.50)	4.79(0.36)	
β_{601}	-4.64(0.17)	-4.44(0.58)	-2.69 (Inf)	-4.65(0.19)	
β_{701}	4.60(0.33)	4.53 (0.45)	4.13(0.87)	4.50(0.38)	
β_{801}	-4.24(0.43)	-4.26(0.58)	-3.85(0.87)	-3.82(1.16)	
β_{901}	4.86(0.15)	$5.01 \ (0.07)$	3.96(1.29)	4.80(0.17)	

Monte Carlo standard errors show high values for many of the coefficients with respect to all the Bayesian regressions. This supports the convergence problems highlighted by the Geweke diagnostic.

Graphical representation of the estimate vs iteration

The trace plots of a set of selected coefficients are reported. Specifically, the estimates of β_1 , β_{401} and β_{801} across iterations are shown for all the regression models. Actual values of the three displayed coefficients are: $\beta_1, \beta_{401}, \beta_{801} = -5$ (depicted with a red line in each graph).



2.6. $n=500, p=1000, \beta^{\star} = \{1, 5, 10\}$

Actual values of the coefficients are: $\beta_1, \beta_{41}, \beta_2, \beta_3, \beta_4 = 1, \beta_5, \beta_6, \beta_7, \beta_8 = 5$ and $\beta_9, \beta_{10} = 10$.

Geweke Diagnostic

G statistic represents the Z-score for a test of equality of means between the first and last parts of the chain (H_0) . A G statistic less than $|z_{\alpha/2}|$, where α is the specific significance level for the test, indicates that data support H_0 , so samples are drawn from the stationary distribution of the chain.

G statistics						
	BL	HS	NG	BR		
β_1	-28.49	-1.62	-53.33	-3.42		
β_2	-4.55	-18.37	-18.83	-45.45		
β_3	-1.76	-5.98	-41.94	-1.79		
β_4	-3.22	-14.23	-41.52	-7.49		
β_5	-1.29	-1.24	-230.25	-168.59		
β_6	-1.54	-2.16	-1.66	-1.23		
β_7	-0.70	-1.98	-0.92	-0.66		
β_8	-1.67	-5.57	-1.45	-1.79		
β_9	-1.38	-2.42	-0.94	-1.20		
β_{10}	1.30	-0.60	-1.45	0.36		

Fixing a significance level $\alpha = 0.01$, we can see that all the approaches generally fail in reaching the convergence for those coefficients characterized by small actual values (from β_1 to β_4).

Monte Carlo standard errors

Estimates of the actual non-zero coefficients are reported with their associated standard error (in parenthesis).

Estimates (standard error)					
	BL	$_{ m HS}$	NG	BR	
β_1	0.63(0.34)	0.15(0.04)	0.50(0.48)	0.72(0.29)	
β_2	0.54(0.20)	$0.19 \ (0.06)$	0.39(0.19)	$0.46\ (0.35)$	
β_3	0.83(0.24)	$0.25 \ (0.12)$	$0.73 \ (0.53)$	0.88(0.22)	
β_4	$0.60 \ (0.25)$	$0.35 \ (0.15)$	0.54(0.38)	0.50(0.24)	
β_5	4.43(0.66)	4.22(0.93)	3.27(2.51)	$3.41 \ (2.19)$	
β_6	4.69(0.23)	4.45(0.15)	4.71(0.19)	4.47(0.42)	
β_7	4.59(0.40)	$0.34 \ (0.18)$	4.48(0.63)	4.70(0.30)	
β_8	4.12(1.22)	$3.83\ (2.03)$	4.16(1.12)	4.03(1.59)	
β_9	8.08(2.13)	$9.91 \ (0.07)$	9.28(0.48)	9.17(0.43)	
β_{10}	$10.37 \ (0.05)$	$10.32 \ (0.59)$	8.56(2.34)	$10.17 \ (0.04)$	

The standard errors are quite similar for all the Bayesian regressions. We notice that the Horseshoe regression presents the lowest standard errors for almost all the estimated coefficients.

Graphical representation of the estimate vs iteration

The trace plots of a set of selected coefficients are reported. Specifically, the estimates of β_1 , β_5 and β_9 across iterations are shown for all the regression models. Actual values of the three displayed coefficients are: $\beta_1 = 1$, $\beta_5 = 5$ and $\beta_9 = 10$ (depicted with a red line in each graph).



2.7. A note on convergence

Having identified some failures to converge, we decided to run the previous simulations increasing the number of MCMC iterations to T = 5000. We then observe that all the analyzed scenarios reach the convergence. The estimates of the actual non-zero coefficients improve their value for some scenarios reducing the Monte Carlo standard errors, even if the approaches still have some difficulties in estimating the correct value of the actual non-zero coefficient for the parameters in the simulations. So convergence doesn't guarantee a better parameter estimation.

3. Increasing the number of MCMC iterations

We develop a set of additional simulations in which the MCMC number of iterations T is increased to evaluate the convergence behavior. Specifically, we consider n = 100 and an increasing number of predictors $p = \{200, 2000, 5000\}$. All simulations refer to $\beta^* = \{1, 5, 10\}$. With respect to the baseline set-up where the number of MCMC iterations equals to T = 1000 with thinning of the chain t = 100 for all the simulations, we increase the number of T to 5000, 10000 and 15000 in correspondence of p equal to 200, 2000 and 5000 respectively. Thinning of the chain t is also increased to 1000 for all simulations. This increase leads to improve the convergence also for those scenarios characterized by low values of n/p ratio. Some difficulties are again detected for $\beta^* = 1$, but for $\beta^* = 5$ and 10 almost all the approaches reach the convergence in particular when the Horseshoe and the Normal-Gamma priors are adopted. The results of these

simulations are provided in the follow.

3.1. p = 200, T = 5000 and t = 1000

Geweke Diagnostic

G statistics								
	BL	HS	NG	BR				
β_1	0.45	0.18	0.53	1.04				
β_2	-0.57	0.63	0.36	3.64				
β_3	1.25	0.06	1.30	0.63				
β_4	-0.55	0.46	0.50	1.12				
β_5	-0.03	-0.38	0.66	0.34				
β_6	-1.11	0.98	-0.77	0.19				
β_7	-0.29	0.31	-0.12	0.99				
β_8	-0.27	0.33	-0.77	1.82				
β_9	-0.68	-1.93	-2.09	0.48				
β_{10}	-0.54	-1.25	-0.31	1.42				

Fixing a significance level $\alpha = 0.01$, all the coefficients reach the convergence.

Monte Carlo standard errors

	Estimates (standard error)									
	BL	HS	NG	BR						
β_1	0.57(0.01)	0.54(0.01)	0.61 (0.01)	0.46(0.01)						
β_2	0.91 (0.01)	0.75~(0.00)	0.89(0.01)	0.91 (0.01)						
β_3	$0.21 \ (0.01)$	0.09(0.00)	$0.16\ (0.00)$	0.20(0.01)						
β_4	0.48(0.01)	$0.51 \ (0.01)$	0.47(0.01)	0.35(0.01)						
β_5	4.34(0.01)	4.61(0.00)	4.50(0.01)	4.21(0.02)						
β_6	4.89(0.01)	5.15(0.00)	5.00(0.00)	4.74(0.01)						
β_7	4.11 (0.01)	4.39(0.00)	4.22(0.01)	3.98(0.01)						
β_8	4.27(0.01)	4.64(0.00)	4.40(0.01)	4.16(0.01)						
β_9	9.12(0.01)	9.46(0.01)	9.26(0.01)	8.63(0.01)						
β_{10}	9.34(0.01)	9.78(0.01)	$9.55\ (0.01)$	9.01 (0.01)						

All the standard errors are very small and almost all the estimated coefficients are in line with their actual values, in particular from β_5 to β_{10} .

Graphical	representation	of	the	estimate	vs	iteration
-----------	----------------	----	-----	----------	----	-----------



3.2. p = 2000, T = 10000 and t = 1000

 $Geweke\ Diagnostic$

G statistics								
	BL	HS	NG	BR				
β_1	-1.05	0.14	0.43	-0.11				
β_2	0.99	-0.17	0.97	0.60				
β_3	-1.09	2.54	0.67	-0.29				
β_4	0.13	-0.65	-0.72	-2.07				
β_5	-0.96	-0.67	-0.91	1.38				
β_6	-0.45	-1.01	1.35	0.15				
β_7	-0.84	-0.83	-1.12	1.12				
β_8	-0.99	-0.75	-0.63	0.73				
β_9	-0.81	-0.68	-1.66	-0.55				
β_{10}	-0.69	-0.99	-0.86	1.56				

Fixing a significance level $\alpha = 0.01$, all the coefficients reach the convergence.

Monte Carlo standard errors

	Estimates (standard error)								
	BL	$_{ m HS}$	NG	BR					
β_1	0.03(0.01)	$0.04 \ (0.01)$	0.03(0.01)	0.00(0.00)					
β_2	0.02(0.00)	$0.02 \ (0.01)$	0.02(0.01)	0.02(0.00)					
β_3	0.18(0.03)	$0.15 \ (0.05)$	0.15(0.03)	0.34(0.04)					
β_4	0.11(0.02)	0.09(0.04)	0.11(0.02)	0.13(0.02)					
β_5	4.11(0.09)	4.32(0.06)	4.31(0.04)	3.01(0.09)					
β_6	3.99(0.01)	4.14(0.06)	4.13(0.04)	3.76(0.03)					
β_7	4.59(0.10)	5.02(0.07)	4.88(0.05)	3.46(0.07)					
β_8	4.80(0.05)	5.01(0.09)	4.88(0.06)	4.44(0.02)					
β_9	9.09(0.17)	9.38(0.07)	9.33(0.02)	7.61(0.03)					
β_{10}	9.49(0.05)	10.02(0.09)	9.74(0.06)	6.70(0.03)					

Almost all the standard errors are very small and the estimated coefficients from β_5 to β_{10} are in line with their actual values.



3.3. p = 5000, T = 15000 and t = 1000

 $Geweke\ Diagnostic$

G statistics								
	BL	HS	NG	BR				
β_1	-0.39	-1.69	0.38	-1.62				
β_2	-1.83	-1.38	-0.37	-0.37				
β_3	-1.28	-1.89	0.16	0.76				
β_4	-0.82	-1.05	-0.04	-2.37				
β_5	-1.78	-3.28	-1.48	-1.42				
β_6	-2.58	-1.45	-2.32	0.90				
β_7	-0.94	-0.53	0.15	-1.65				
β_8	-2.13	-2.08	-0.13	-0.85				
β_9	-1.55	-1.24	0.31	-90.18				
β_{10}	-1.15	-0.85	-0.75	-1.66				

Fixing a significance level $\alpha = 0.01$, all the coefficients reach the convergence except for β_5 in the Horseshoe regression and β_9 in the Bayesian Ridge.

Monte Carlo standard errors

	Estimates (standard error)									
	BL	HS	NG	BR						
β_1	0.03(0.01)	$0.03\ (0.03)$	$0.01 \ (0.00)$	0.01 (0.00)						
β_2	0.02(0.01)	0.01(0.01)	$0.03 \ (0.01)$	$0.00 \ (0.00)$						
β_3	0.01 (0.00)	0.02(0.01)	$0.01 \ (0.00)$	0.00(0.00)						
β_4	0.06(0.02)	0.07(0.04)	0.13(0.03)	0.00(0.00)						
β_5	3.76(0.20)	3.82(0.41)	3.81(0.29)	2.32(0.42)						
β_6	3.52(0.34)	4.05(0.30)	4.06(0.05)	0.64(0.18)						
β_7	4.85(0.15)	5.06(0.10)	5.05(0.03)	4.27(0.04)						
β_8	3.60(0.31)	4.02(0.40)	4.22(0.03)	0.79(0.26)						
β_9	9.25(0.03)	9.33(0.60)	9.49(0.04)	5.59(0.78)						
β_{10}	8.94 (0.49)	$9.67\ (0.35)$	9.62(0.05)	5.46(0.12)						

The standard errors are small for many of the estimated coefficients and the Normal-Gamma regression presents almost always the smallest values of standard errors. The estimates of the coefficients with high actual values (specifically, from β_5 to β_{10}) are generally close to their actual values.

Graphical	representation	of	the	estimate	vs	iteration
-----------	----------------	----	-----	----------	----	-----------



3.4. Computational cost

To quantify the computational cost of the simulations, we derive the CPU and the total elapsed time spent (in seconds) when T = 1000 and t = 100 (the baseline values for T and t used in the paper) and when T = 5000, 10000, 15000 and t = 1000 (the increased values of T and t used in the supplemental simulations conducted to reach convergence). Time refers to only one MCMC run, and from the results reported in Table 5 we can see how the procedures increase their computational time making difficult the development of a high numerosity of scenarios for the comparison.

[Table 5 about here.]

4. Convergence on real data example

We provide the results of the convergence on the real dataset described in Section 4 of the paper. We recall that these data are characterized by a number of variables (fragments) p = 4059 and a number of observations (molecules) n = 140 sampled by the whole set of N = 1704 active molecules. In this analysis we set for T = 15000 and t = 1000.

Geweke Diagnostic

G statistic represents the Z-score for a test of equality of means between the first and last parts of the chain (H_0) . A G statistic less than $|z_{\alpha/2}|$, where α is the specific

	G statistics									
	BL	HS	NG	BR						
β_{F14}	1.19	0.76	1.38	0.77						
β_{F24}			-2.81							
β_{F80}	1.37	0.76	1.23	1.36						
β_{F119}		-1.27								
β_{F144}		2.95		-0.75						
β_{F177}		0.55	-2.06							
β_{F219}	-0.31			0.41						
β_{F222}	-1.63		-7.09							
β_{F253}				-1.10						

significance level for the test, indicates that data support H_0 , so samples are drawn from the stationary distribution of the chain.

Fixing a significance level $\alpha = 0.01$, all the coefficients reach the convergence except for β_{F144} in the Horseshoe regression, β_{F24} and β_{F222} in the Normal-Gamma regression.

Monte Carlo standard errors

Estimates of the actual non-zero coefficients are reported with their associated standard error (in parenthesis).

	Estimates (standard error)										
	BL	HS	NG	BR							
β_{F14}	-2.12(0.02)	-2.12(0.04)	-2.14(0.01)	-2.13(0.03)							
β_{F24}			-0.50(0.12)								
β_{F80}	-1.93(0.13)	-2.03(0.02)	-2.02(0.01)	-1.75(0.12)							
β_{F119}		0.85(0.10)									
β_{F144}		-0.89(0.17)		-0.82(0.18)							
β_{F177}		-0.67(0.08)	-0.74(0.10)								
β_{F219}	-0.77(0.14)			-0.75 (0.10)							
β_{F222}	-0.92(0.14)		-1.03(0.23)								
β_{F253}				0.77(0.15)							

Almost all the standard errors are very small and the estimates of the common coefficients are very similar across all the Bayesian approaches.

Graphical representation of the estimate vs iteration

The trace plots of the two coefficients selected by all the approaches are reported. Specifically, the estimates of β_{F14} and β_{F80} across iterations are shown for all the regression models.



References

- Griffin JE, Brown PJ. Inference with normal-gamma prior distributions in regression problems. Bayesian Analysis. 2010 03;5(1):171–188.
- [2] Geweke J. Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments. In: Bernardo JM, Berger JO, Dawid AP, Smith AFM, editors. Bayesian Statistics 4; Oxford University Press, 1992. p. 169–193.
- [3] Jones GL, Haran M, Caffo BF, Neat R. Fixed-Width Output Analysis for Markov Chain Monte Carlo. Journal of the American Statistical Association. 2006; 101(476):1537–1547.

Table 1. Predictive Mean Square Error for regression with $\beta^* = \{0.1, 0.5\}$. BL=Bayesian Lasso; HS= Horseshoe regression; NG= Normal-Gamma regression; BR= Bayesian Ridge (standard errors of MonteCarlo runs in parentheses).

			$\beta^{\star} = 0.1$		
\overline{n}	p	BL	HS	NG	BR
100	200	1.06(0.03)	1.05(0.04)	1.06(0.03)	1.45(0.14)
100	5000	1.02(0.03)	1.02(0.04)	1.02(0.03)	$1.21 \ (0.06)$
			$\beta^{\star} = 0.5$		
		BL	HS	NG	BR
100	200	1.25(0.04)	1.27(0.08)	1.25(0.05)	1.58(0.12)
100	5000	1.22(0.03)	1.22(0.05)	1.22(0.03)	1.43(0.05)

Table 2. Total Number of Discoveries in regression simulations with $k^{\star} = 10$ non-zero coefficients β , BL=Bayesian Lasso; HS= Horseshoe regression; NG= Normal-Gamma regression; BR= Bayesian Ridge. Median values of TND achieved in the 50 MonteCarlo runs are reported.

			$\beta^{\star} = 0.1$			$\beta^{\star} = 0.5$			
\overline{n}	p	BL	HS	NG	BR	BL	HS	NG	BR
100	200	36	30	35	54	39	34	37	55
100	5000	7	5	6	20	6	6	7	20

Table 3. False Discovery Rate for regression with $\beta^* = \{0.1, 0.5\}$. BL=Bayesian Lasso; HS= Horseshoe regression; NG= Normal-Gamma regression; BR= Bayesian Ridge (standard errors of MonteCarlo runs in parentheses).

$\beta^{\star} = 0.1$									
\overline{n}	p	BL	HS	NG	BR				
100	200	$0.93 \ (0.03)$	0.93(0.04)	0.93(0.04)	$0.94 \ (0.02)$				
100	5000	1.00(0.00)	0.99~(0.05)	1.00(0.00)	1.00(0.00)				
$\beta^{\star} = 0.5$									
		BL	HS	NG	BR				
100	200	0.88(0.04)	0.86(0.05)	0.87(0.05)	0.90(0.03)				
100	5000	1.00(0.01)	0.98(0.08)	0.99(0.04)	0.99(0.01)				

$\beta^* = 0.1$									
n	p BL		HS	NG	BR				
100	200	0.26(0.12)	0.21(0.12)	0.25(0.13)	0.33(0.13)				
100	5000	0.00(0.00)	$0.01 \ (0.02)$	$0.00\ (0.00)$	$0.01 \ (0.03)$				
$\beta^{\star} = 0.5$									
		BL	HS	NG	BR				
100	200	$0.51 \ (0.15)$	0.47(0.16)	0.48(0.19)	0.56(0.16)				
100	5000	0.01(0.01)	0.01 (0.02)	0.01(0.02)	0.02(0.04)				

Table 4. Sensitivity for regression with $\beta^* = \{0.1, 0.5\}$. BL=Bayesian Lasso; HS= Horseshoe regression; NG= Normal-Gamma regression; BR= Bayesian Ridge (standard errors of MonteCarlo runs in parentheses).

Table 5. Computational cost in terms of CPU and total elapsed time spent (in seconds) of one MCMC run. All the scenarios refer to n = 100 and $\beta^* = \{1, 5, 10\}$. BL=Bayesian Lasso; HS= Horseshoe regression; NG= Normal-Gamma regression; BR= Bayesian Ridge.

	CPU time				Elapsed time			
~ .				Diapsed time				
Scenario	BL	HS	NG	BR	BL	HS	NG	BR
p = 200, T = 1000 and t = 100	0.250	0.269	0.540	0.352	143.044	174.137	177.697	86.438
p = 200, T = 5000 and $t = 1000$	15.680	8.898	8.844	2.854	7025.436	9234.502	8668.242	2765.421
p = 2000, T = 1000 and t = 100	0.420	0.494	0.704	0.485	251.445	251.788	260.201	249.553
p = 2000, T = 10000 and $t = 1000$	15.671	15.997	17.365	16.747	24302.849	24375.102	25192.927	24195.071
p = 5000, T = 1000 and t = 100	0.474	0.438	0.317	0.229	324.701	251.997	254.521	244.805
p = 5000, T = 15000 and t = 1000	25.398	28.438	25.019	28.551	36570.723	36743.195	37903.543	36884.321