

The Illusion of the *Illusion of Sparsity*

BRUNO FAVA

PhD Student

Northwestern University, Illinois, USA

HEDIBERT F. LOPES

Professor of Statistics and Econometrics

Head of the Center of Statistics, Data Science and Decision

INSPER, São Paulo, Brazil

January 2021

Manuscript and slides found at hedibert.org

Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

- Ridge and lasso regressions

- Spike and slab model (or SMN model)

- SSVS and scaled SSVS priors

- R package `Bayeslm`

Recalling GLP main remarks

- Inclusion and tail probabilities

Our experiments

- I. Adding meaningless variables

- II. Fatter tails via Student's t

- III. More simulations

Final remarks

Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

- Ridge and lasso regressions

- Spike and slab model (or SMN model)

- SSVS and scaled SSVS priors

- R package `Bayeslm`

Recalling GLP main remarks

- Inclusion and tail probabilities

Our experiments

- I. Adding meaningless variables

- II. Fatter tails via Student's t

- III. More simulations

Final remarks

Sparsity in macro, micro and finance

We revisit Giannone, Lenza and Primiceri's (GLP) *Economic predictions with big data: the illusion of sparsity*, whose abstract says ¹:

1. We compare **sparse** and **dense** representations of predictive models in **macroeconomics, microeconomics and finance**.
2. To deal with a large number of possible predictors, we specify a **prior** that allows for both **variable selection** and **shrinkage**.
3. The posterior distribution does not typically concentrate on a **single sparse model**, but on a wide set of models that often include **many predictors**.

They conclude:

- ▶ No reason **predictive models** should include only a handful of predictors.
- ▶ **Low-dimensional models** justified only with **strong statistical evidence**.

¹<https://faculty.wcas.northwestern.edu/~gep575/illusion4-2.pdf>

GLP datasets

| | Dependent variable | Possible predictors | Sample |
|------------------|---|--|---|
| Macro 1 | Monthly growth rate of US industrial production | 130 lagged macroeconomic indicators | 659 monthly time-series observations, from February 1960 to December 2014 |
| Macro 2 | Average growth rate of GDP over the sample 1960-1985 | 60 socio-economic, institutional and geographical characteristics, measured at pre-1960s value | 90 cross-sectional country observations |
| Finance 1 | US equity premium (S&P 500) | 16 lagged financial and macroeconomic indicators | 68 annual time-series observations, from 1948 to 2015 |
| Finance 2 | Stock returns of US firms | 144 dummies classifying stock as very low, low, high or very high in terms of 36 lagged characteristics | 1400k panel observations for an average of 2250 stocks over a span of 624 months, from January 1963 to May 2014 |
| Micro 1 | Per-capita crime (murder) rates | Effective abortion rate and 284 controls including possible covariate of crime and their transformations | 576 panel observations for 48 US states over a span of 144 months, from 1986 to 1997 |
| Micro 2 | Number of pro-plaintiff eminent domain decisions in a specific circuit and in a specific year | Characteristics of judicial panels capturing aspects related to gender, race, religion, political affiliation, education and professional history of the judges, together with some interactions among the latter, for a total of 138 regressors | 312 panel observations for 12 circuits over a span of 26 years, from 1979 to 2004 |

TABLE 1. Description of the datasets.

Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

- Ridge and lasso regressions

- Spike and slab model (or SMN model)

- SSVS and scaled SSVS priors

- R package `Bayeslm`

Recalling GLP main remarks

- Inclusion and tail probabilities

Our experiments

- I. Adding meaningless variables

- II. Fatter tails via Student's t

- III. More simulations

Final remarks

GLP spike-and-slab prior

They consider a standard Gaussian linear model:

$$y_t = \beta_1 x_{t1} + \cdots + \beta_k x_{tk} + \epsilon_t, \quad \epsilon_t \sim N(0, \sigma^2).$$

The prior for β_i accommodates sparsity and/or shrinkage:

$$\beta_i | \sigma^2, \gamma^2, q \sim \begin{cases} N(0, \sigma^2 \gamma^2) & \text{with prob. } q \\ 0 & \text{with prob. } 1 - q \end{cases} \quad i = 1, \dots, k.$$

q governs the **degree of sparsity**.

γ^2 governs the **degree of shrinkage**.

Alternative representation

An alternative way of writing the prior for β_i is

$$\beta_i | \sigma^2, \gamma^2, \mathbf{q}, \nu_i \sim N(0, \sigma^2 \gamma^2 \nu_i) \quad \text{and} \quad \nu_i \sim \text{Ber}(q).$$

The **Bayesian lasso**, **horseshoe** and **elastic net** methods can instead be obtained by replacing the Bernoulli for ν_i with an **exponential**, a **half-Cauchy**, or a **transformation of a truncated Gamma**, respectively.

None of these priors admit a truly sparse representation with positive probability.

We will get back to this representation soon!

Hyperprior of (q, γ^2)

Instead of setting a hyperprior for (q, γ^2) , GLP defined a prior for the pair (q, R^2) , where

$$R^2(\gamma^2, q) \equiv \frac{qk\gamma^2}{qk\gamma^2 + 1},$$

is the coefficient of determination.

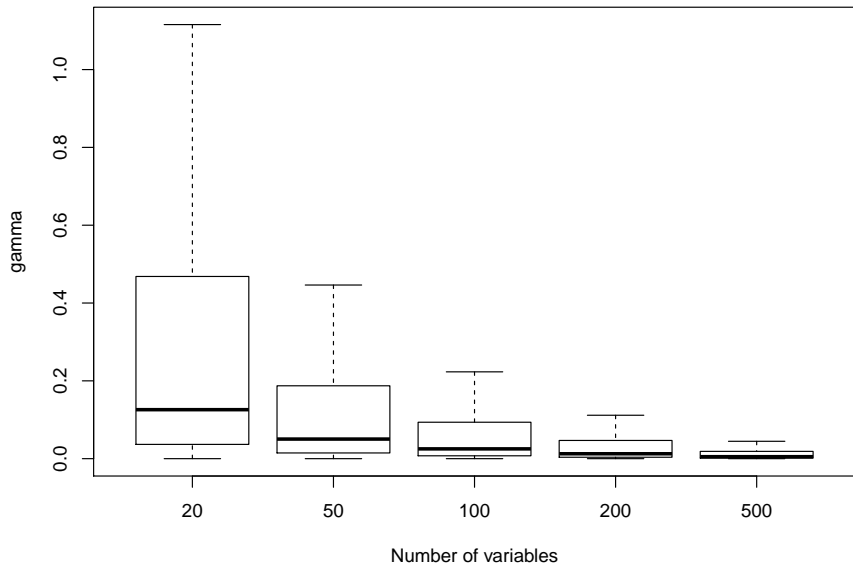
The hyperprior distributions are:

$$q \sim \text{Beta}(1, 1)$$

and

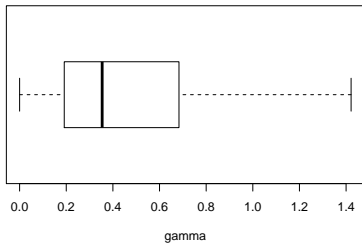
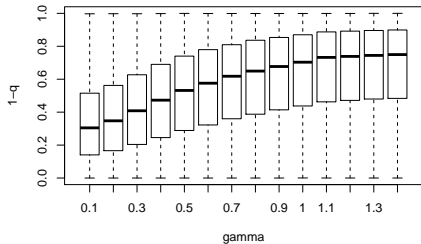
$$R^2 \sim \text{Beta}(1, 1)$$

Marginal prior of γ : $p(\gamma|k)$

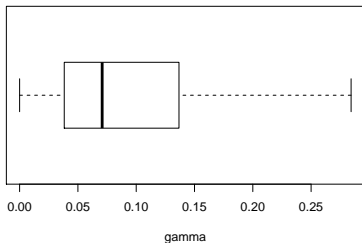
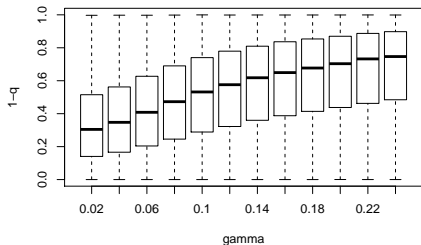


$p(1 - q|\gamma)$ and $p(\gamma)$

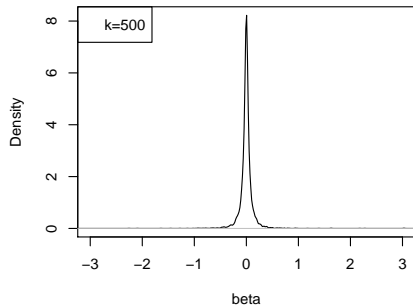
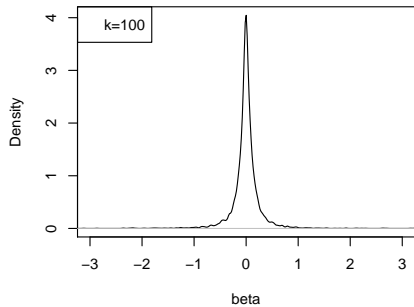
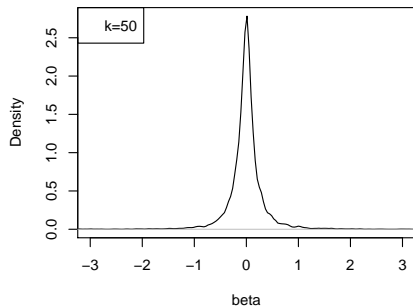
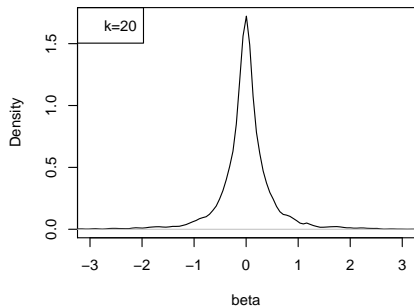
Pr(q|gamma,k=20)



Pr(q|gamma,k=500)



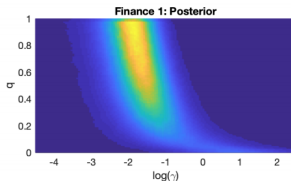
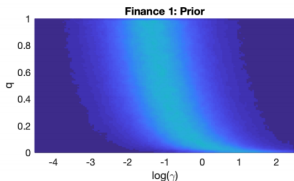
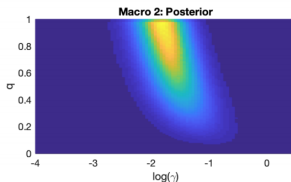
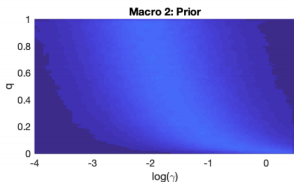
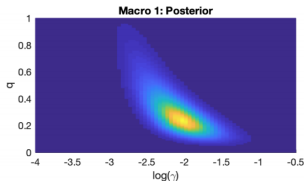
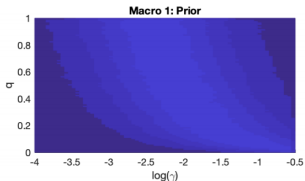
$$p(\beta|k, \sigma = 1)$$



GLP: $p(q, \gamma | data)$ - Macro 1, Macro 2 and Finance1

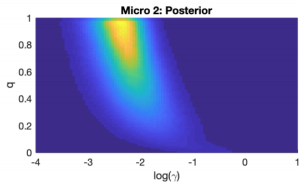
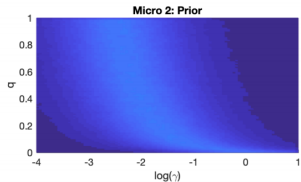
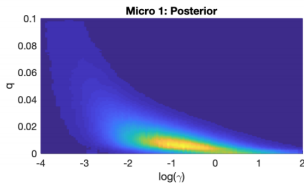
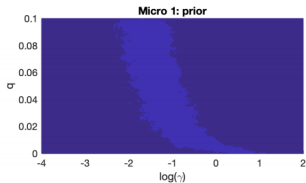
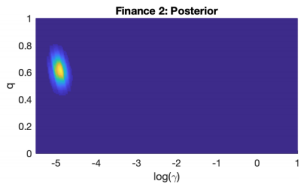
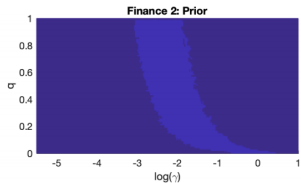
Negative correlation:

higher the probability of inclusion, lower prior variance of a non-zero coefficient.

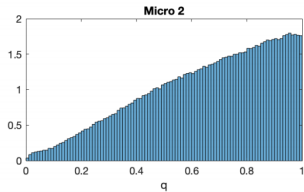
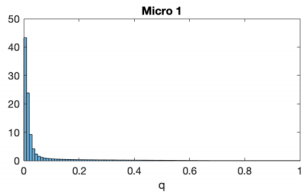
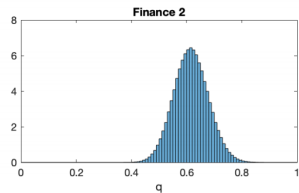
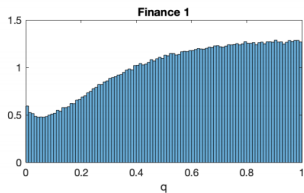
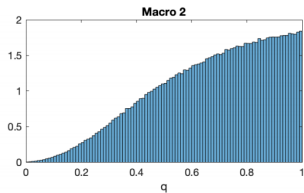
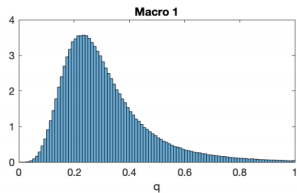


GLP: $p(q, \gamma | data)$ - Finance 2, Micro 1 and Micro 2

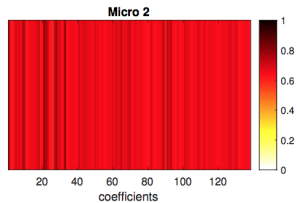
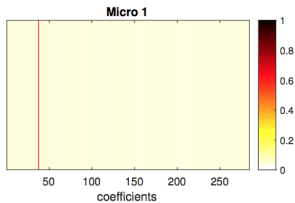
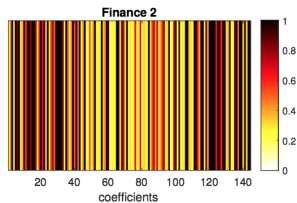
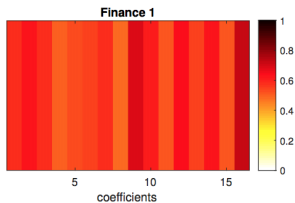
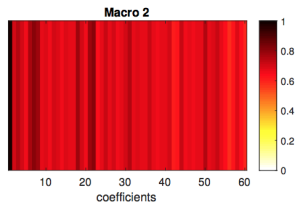
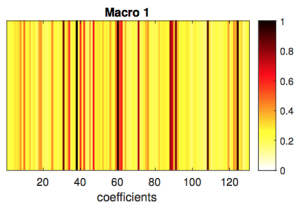
Micro 1 is the only case where $p(q|data)$ is concentrated around very low values.



GLP: $p(q|data)$



GLP: Probability of inclusion of each predictor



Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

- Ridge and lasso regressions

- Spike and slab model (or SMN model)

- SSVS and scaled SSVS priors

- R package `Bayeslm`

Recalling GLP main remarks

- Inclusion and tail probabilities

Our experiments

- I. Adding meaningless variables

- II. Fatter tails via Student's t

- III. More simulations

Final remarks

Our contribution

- ▶ We analyze the posterior distribution of the included coefficients of the linear model. This was not explored by GLP
- ▶ We add meaningless predictors and observe correct exclusion only in a subset of the simulated data sets.
- ▶ We consider a Student's t prior $\beta_i | \beta \neq 0$.
 - ▶ More restrictive in selecting possible predictors.
- ▶ We show, via simulations, that their prior incorrectly induces shrinkage.

Overall conclusion: Their Spike-and-Slab approach does not seem to be robust, leading to the illusion that sparsity is nonexistent, when it might in fact exist. Therefore, *the illusion of the illusion :o)*

Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

Ridge and lasso regressions

Spike and slab model (or SMN model)

SSVS and scaled SSVS priors

R package `Bayeslm`

Recalling GLP main remarks

Inclusion and tail probabilities

Our experiments

I. Adding meaningless variables

II. Fatter tails via Student's t

III. More simulations

Final remarks

Ridge and lasso regressions

Recall the standard Gaussian linear model,

$$y_t = \beta_1 x_{1t} + \beta_2 x_{2t} + \cdots + \beta_k x_{kt} + \nu_t,$$

where $RSS = (y - X\beta)'(y - X\beta)$ is the residual sum of squares.

- ▶ *Ridge regression* Hoerl and Kennard [1970] - ℓ_2 penalty on β :

$$\hat{\beta}_{ridge} = \arg \min_{\beta} \left\{ RSS + \lambda_r^2 \sum_{j=1}^k \beta_j^2 \right\}, \quad \lambda_r^2 \geq 0,$$

leading to $\hat{\beta}_{ridge} = (X'X + \lambda_r^2 I_k)^{-1} X'y$.

- ▶ *Lasso regression* Tibshirani [1996] - ℓ_1 penalty on β :

$$\hat{\beta}_{lasso} = \arg \min_{\beta} \left\{ RSS + \lambda_l \sum_{j=1}^k |\beta_j| \right\}, \quad \lambda_l \geq 0,$$

which can be solved by a *coordinate gradient descent* algorithm.

Ridge and lasso estimates are posterior modes!

The posterior mode or the maximum a posteriori (MAP) is given by

$$\tilde{\beta}_{\text{mode}} = \arg \min_{\beta} \{-2 \log p(y|\beta) - 2 \log p(\beta)\}$$

The $\hat{\beta}_{\text{ridge}}$ estimate equals the posterior mode of the normal linear model with

$$p(\beta_j) \propto \exp\{-0.5\lambda_r^2\beta_j^2\},$$

which is a **Gaussian distribution** with location 0 and scale $1/\lambda_r^2$, $N(0, 1/\lambda_r^2)$. The mean is 0, the variance is $1/\lambda_r^2$ and the excess kurtosis is 0.

The $\hat{\beta}_{\text{lasso}}$ estimate equals the posterior mode of the normal linear model with

$$p(\beta_j) \propto \exp\{-0.5\lambda_l|\beta_j|\},$$

which is a **Laplace distribution** with location 0 and scale $2/\lambda_l$, $\text{Laplace}(0, 2/\lambda_l)$. The mean is 0, the variance is $8/\lambda_l^2$ and excess kurtosis is 3.

Spike and slab model (or scale mixture of normals)

Ishwaran and Rao [2005] define a **spike and slab model** as a Bayesian model specified by the following prior hierarchy:

$$\begin{aligned}(y_t | x_t, \beta, \sigma^2) &\sim N(x_t' \beta, \sigma^2), & t = 1, \dots, n \\ (\beta | \psi) &\sim N(0, \text{diag}(\psi)) \\ \psi &\sim \pi(d\psi) \\ \sigma^2 &\sim \mu(d\sigma^2)\end{aligned}$$

They go to say that

“Lempers [1988] and Mitchell and Beauchamp [1988] were among the earliest to pioneer the spike and slab method. The expression ‘spike and slab’ referred to the prior for β used in their hierarchical formulation.”

Spike and slab model (or scale mixture of normals model)

Regularization and variable selection are done by assuming independent prior distributions from the SMN class to each coefficient β_j :

$$\beta_j | \psi_j \sim N(0, \psi_j) \quad \text{and} \quad \psi_j \sim p(\psi_j)$$

so

$$p(\beta_j) = \int p(\beta_j | \psi_j) p(\psi_j) d\psi_j.$$

| Mixing density $p(\psi_j)$ | Marginal density $p(\beta_j)$ | $V(\beta_j)$ | Ex.kurtosis(β_j) |
|----------------------------|---------------------------------------|-----------------------|--------------------------|
| $\psi_j = 1/\lambda_r^2$ | $N(0, 1/\lambda_r^2)$ - (ridge) | $1/\lambda_r^2$ | 0 |
| $IG(\eta/2, \eta\tau^2/2)$ | $t_\eta(0, \tau^2)$ | $\eta/(\eta-2)\tau^2$ | $6/(\eta-4)$ |
| $G(1, \lambda_l^2/8)$ | Laplace(0, $2/\lambda_l$) - (blasso) | $8/\lambda_l^2$ | 3 |
| $G(\zeta, 1/(2\gamma^2))$ | $NG(\zeta, \gamma^2)$ | $2\zeta\gamma^2$ | $3/\zeta$ |

Griffin and Brown [2010] Normal-Gamma prior:

$$p(\beta | \zeta, \gamma^2) = \frac{1}{\sqrt{\pi} 2^{\zeta-1/2} \gamma^{\zeta+1/2} \Gamma(\zeta)} |\beta|^{\zeta-1/2} K_{\zeta-1/2}(|\beta|/\gamma),$$

where K is the modified Bessel function of the 3rd kind.

Illustration

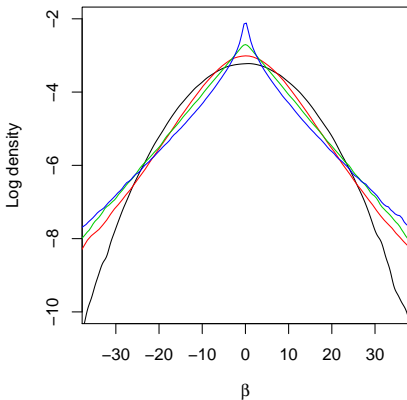
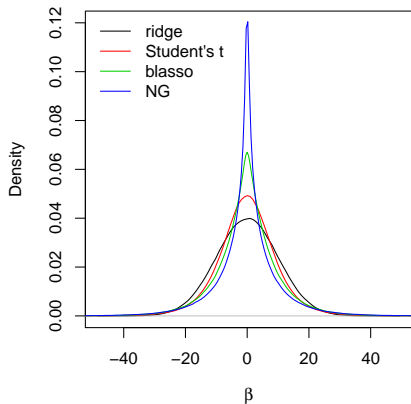
Ridge: $\lambda_r^2 = 0.01 \Rightarrow$ Excess kurtosis=0

Student's t : $\eta = 5, \tau^2 = 60 \Rightarrow$ Excess kurtosis=6

Blasso: $\lambda_l^2 = 0.08 \Rightarrow$ Excess kurtosis=3

NG: $\xi = 0.5, \gamma^2 = 100 \Rightarrow$ Excess kurtosis=6

All variances are equal to 100.



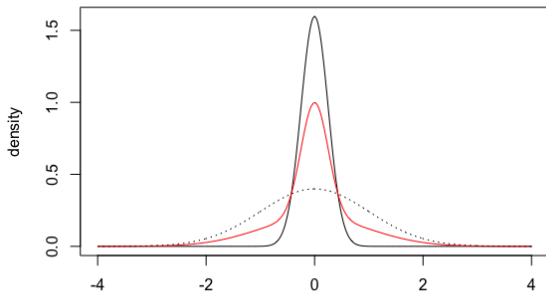
Stochastic search variable selection (SSVS) prior

SSVS George and McCulloch [1993]: For small $\tau > 0$ and $c \gg 1$,

$$\beta | \omega, \tau^2, c^2 \sim (1 - \omega) \underbrace{N(0, \tau^2)}_{\text{spike}} + \omega \underbrace{N(0, c^2 \tau^2)}_{\text{slab}}.$$

SMN representation: $\beta | \psi \sim N(0, \psi)$ and

$$\psi | \omega, \tau^2, c^2 \sim (1 - \omega) \delta_{\tau^2}(\psi) + \omega \delta_{c^2 \tau^2}(\psi)$$

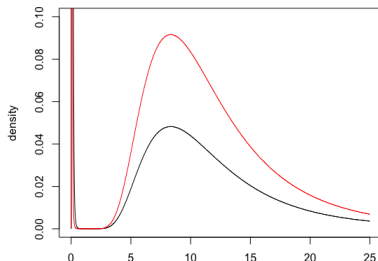


Scaled SSVS prior = normal mixture of IG prior

NMIG prior of Ishwaran and Rao [2005]: For $v_0 \ll v_1$,

$$\begin{aligned}\beta|K, \tau^2 &\sim N(0, K\tau^2), \\ K|\omega, v_0, v_1 &\sim (1 - \omega)\delta_{v_0}(K) + \omega\delta_{v_1}(K), \\ \tau^2 &\sim IG(a_\tau, b_\tau).\end{aligned}\tag{1}$$

- ▶ Large ω implies non-negligible effects.
- ▶ The scale $\psi = K\tau^2 \sim (1 - \omega)IG(a_\tau, v_0b_\tau) + \omega IG(a_\tau, v_1b_\tau)$.
- ▶ $p(\beta)$ is a **two component mixture of scaled Student's t distributions**.



R package Bayeslm

Bayeslm was written by Jingyu He and is based on Hahn, He and Lopes (2019) **Efficient sampling for Gaussian linear regression with arbitrary priors**, *Journal of Computational and Graphical Statistics*, 28, 142-154.

For observation $i = 1, \dots, n = 68$ and predictor $j = 1, \dots, k = 16$, we simulate

$$x_{ij} \sim N(0, 1) \quad \text{and} \quad \varepsilon_i^* \sim N(0, 1)$$

We also fix $\beta_1 = -0.86$, $\beta_2 = 0.64$ and $\beta_3 = 0.89$, while the response variable is:

$$y_i^{(s)} = \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \sigma_\varepsilon^{(s)} \varepsilon_i^*,$$

and $\sigma_\varepsilon^{(s)} = 0.75s$, for $s = 1, 2$.

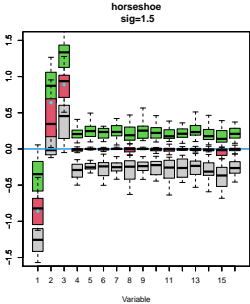
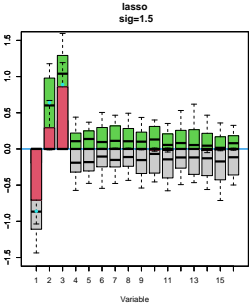
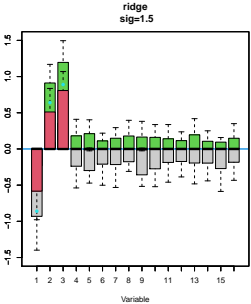
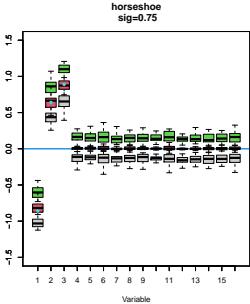
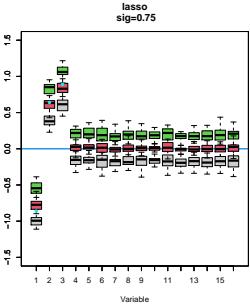
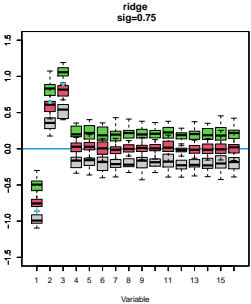
MCMC set-up: $N = 2000$ draws, burnin= 10000 burn-in

Monte Carlo error: $R = 20$ replicates

R script

```
install.packages("bayeslm");library("bayeslm")
n=68;k=16;betas=c(-0.86,0.64,0.89,rep(0,k-3));sigs=c(0.75,1.5)
N=2000;burnin=10000;R=20
qs=c(0.025,0.5,0.975)
J=length(sigs);quants=array(0,c(R,J,3,k,3))
set.seed(54321)
for (r in 1:R){
  for (j in 1:J){
    X = matrix(rnorm(n*k),n,k)
    y = rnorm(n,X%*%betas,sigs[j])
    fit.hs = bayeslm(y,x,prior='horseshoe',N=N,burnin=burnin,icept=FALSE)
    fit.ridge = bayeslm(y,x,prior='ridge',N=N,burnin=burnin,icept=FALSE)
    fit.lasso = bayeslm(y,x,prior='laplace',N=N,burnin=burnin,icept=FALSE)
    quants[r,j,1,,] = t(apply(fit.hs$beta,2,quantile,qs))
    quants[r,j,2,,] = t(apply(fit.ridge$beta,2,quantile,qs))
    quants[r,j,3,,] = t(apply(fit.lasso$beta,2,quantile,qs))
  }
}
method = c("horseshoe","ridge","lasso")
par(mfrow=c(2,3))
for (i in 1:2)
  for (j in c(2,3,1)){
    boxplot(quants[,i,j,,1],names=1:k,ylim=c(-1.5,1.5),outline=FALSE,col=gray(0.8),
            xlab="Variable",main=paste(method[j],"\n sig=",sigs[i],sep=""))
    abline(h=0,col=4,lwd=2)
    for (l in 3:2)
      boxplot(quants[,i,j,,l],names=rep("",k),outline=FALSE,col=1,add=TRUE)
    points(1:3,betas[1:3],col=5,pch=16)
  }
}
```

Ridge, Laplace and horseshoe priors

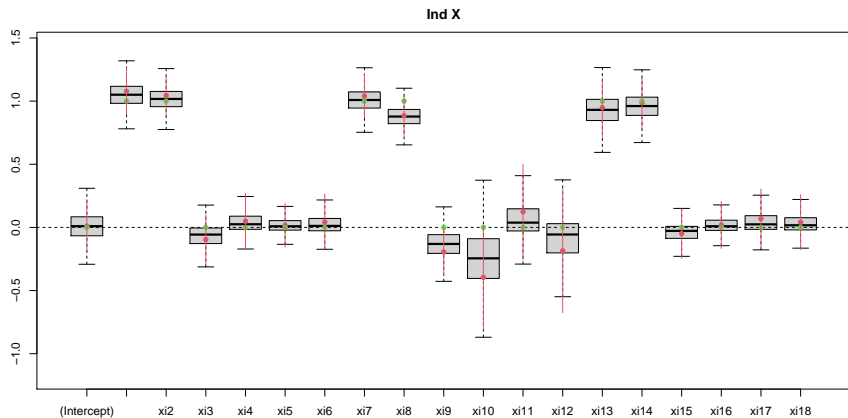


Various designs for X

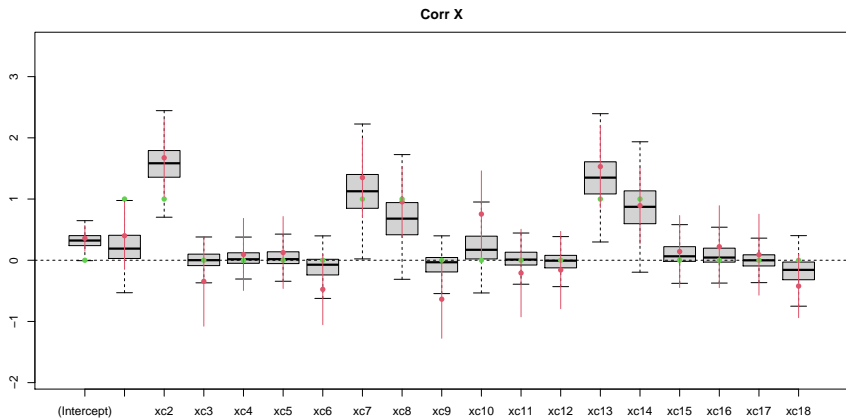
- 1) $\rho(x_i, x_j) = 0.0$ for all i, j
- 2) $\rho(x_i, x_j) = 0.9$ for all i, j
- 3) Two common factors generate the k predictors

```
library(bayeslm)
set.seed(12345)
k      = 18
n      = 100
sig    = 2.0
sig    = 1.0
tau    = 0.1
theta1 = c(rnorm(k/2,1,0.1),rnorm(k/2,0,0.1))
beta   = rep(0,k)
beta[c(1:2,7:8,13:14)] = 1
theta2 = c(rnorm(k/3,0,0.1),rnorm(k/3,0.5,0.1),rnorm(k/3,1,0.1))
theta  = cbind(theta1,theta2)
fac    = matrix(rnorm(2*n),n,2)
xf     = fac%*%t(theta)+rnorm(n*k,0,tau)
Vx     = diag(apply(xf,2,var))
yf     = xf%*%beta+rnorm(n,0,sig)
xi     = matrix(rnorm(n*k),n,k)%*%sqrt(Vx)
yi     = xi%*%beta+rnorm(n,0,sig)
Sigma  = matrix(1,k,k)
rho    = 0.9
for (i in 1:k)
  for (j in 1:k)
    Sigma[i,j] = rho^(abs(i-j))
xc     = matrix(rnorm(n*k),n,k)%*%chol(Sigma)
yc     = xc%*%beta+rnorm(n,0,sig)
```

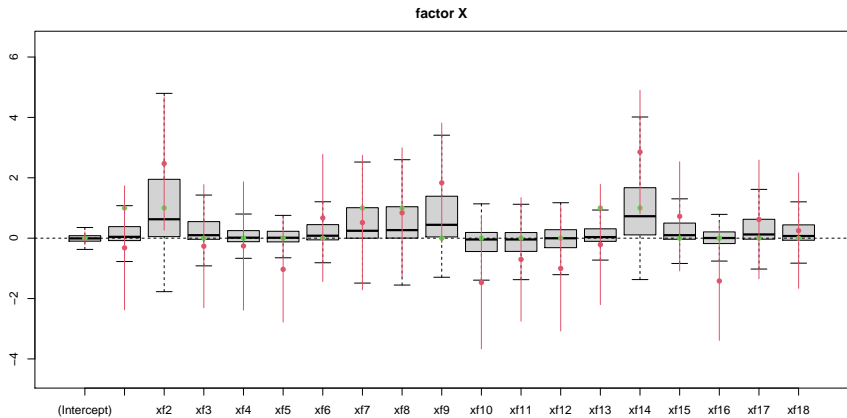
Independent predictors



Correlated predictors



X matrices



A few additional references

Park and Casella (2008) The Bayesian lasso. *JASA*, 103(482), 681-686.

Carvalho, Polson and Scott (2010) The horseshoe estimator for sparse signals. *Biometrika*, 97(2)465-480.

Polson and Scott (2010) Shrink globally, act locally: Sparse Bayesian regularization and prediction, *Bayesian Statistics*, Volume 9, 501–538.

Polson and Scott (2012) Local shrinkage rules, Lévy processes and regularized regression, *JRSS-B*, 74(2), 287-311.

van der Pas, Kleijn and van der Vaart (2014) The horseshoe estimator: Posterior concentration around nearly black vectors. *Electronic Journal of Statistics*, 8, 2585-2618.

Bhattacharya, Pati, Pillai and Dunson (2015) Dirichlet–Laplace priors for optimal shrinkage, *JASA*, 110, 1479–1490.

Makalic and Schmidt (2016) A Simple Sampler for the Horseshoe Estimator. *IEEE Signal Processing Letters*, 23(1), 179-182.

Bhadra, Datta, Polson and Willard (2017) The Horseshoe+ Estimator of Ultra-Sparse Signals, *Bayesian Analysis*, 12(4), 1105–1131.

Rocková and George (2018) The Spike-and-Slab LASSO, *JASA*, 113(521), 431-444.

Hahn, He and Lopes (2019) Efficient sampling for Gaussian linear regression with arbitrary priors, *JCGS*, 28, 142-154.

Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

Ridge and lasso regressions

Spike and slab model (or SMN model)

SSVS and scaled SSVS priors

R package `Bayeslm`

Recalling GLP main remarks

Inclusion and tail probabilities

Our experiments

I. Adding meaningless variables

II. Fatter tails via Student's t

III. More simulations

Final remarks

Recalling GLP main remarks

The conclusion is that a clear pattern of sparsity is found only on the Micro 1 data set, in which only one variable is included most of the times.

For all other data sets, one is incapable of determining which variables should be included, as many have a high estimated probability of inclusion \Rightarrow **dense models**.

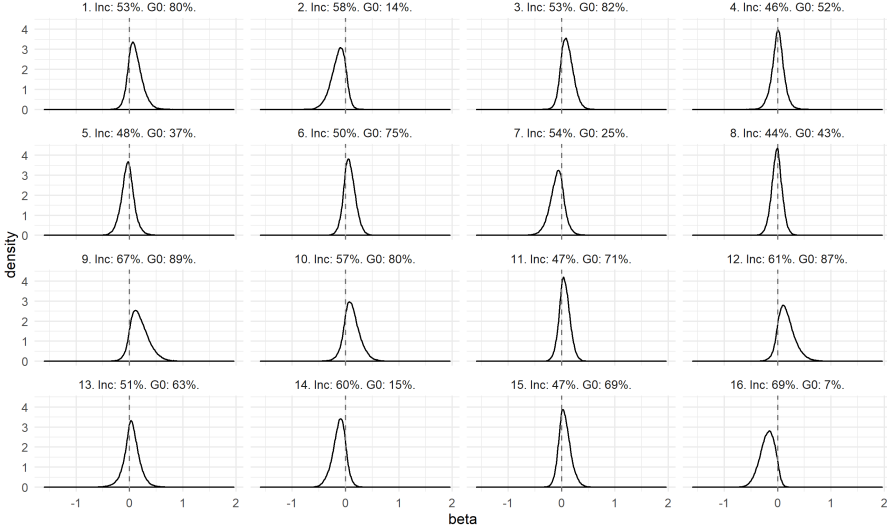
Their conclusion: Sparsity cannot be assumed for any economic data set, unless in the presence of strong statistical evidence, and suggest an "**illusion of sparsity**" when using statistical models that assume (and force) sparsity.

Inclusion and tail probabilities

Finance 1

Inc: Probability of inclusion

G0: Tail (above zero) probability.



An example: $\beta_1, \beta_6, \beta_7$ and β_{13}

The spike-and-slab prior, as defined, seems to be inducing shrinkage by including predictors with a **near-zero** coefficient.

Example: $\beta_1, \beta_6, \beta_7$ and β_{13}

- ▶ Probability of inclusion between 0.5 and 0.54, but also tail probability between 0.2 and 0.4.
- ▶ It could be, for example, that an economist trying to make inference on the regression would very easily exclude variable 1, but keep, for example, variables 6, 7 and 13.

Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

- Ridge and lasso regressions

- Spike and slab model (or SMN model)

- SSVS and scaled SSVS priors

- R package `Bayeslm`

Recalling GLP main remarks

- Inclusion and tail probabilities

Our experiments

- I. Adding meaningless variables

- II. Fatter tails via Student's t

- III. More simulations

Final remarks

I. Adding meaningless variables

We re-run the estimation algorithm for all the five datasets but now include two additional regressors that were completely randomly generated.

Micro 1: 1.6% and 3.9%

Macro 1: 12.2% and 21.1%

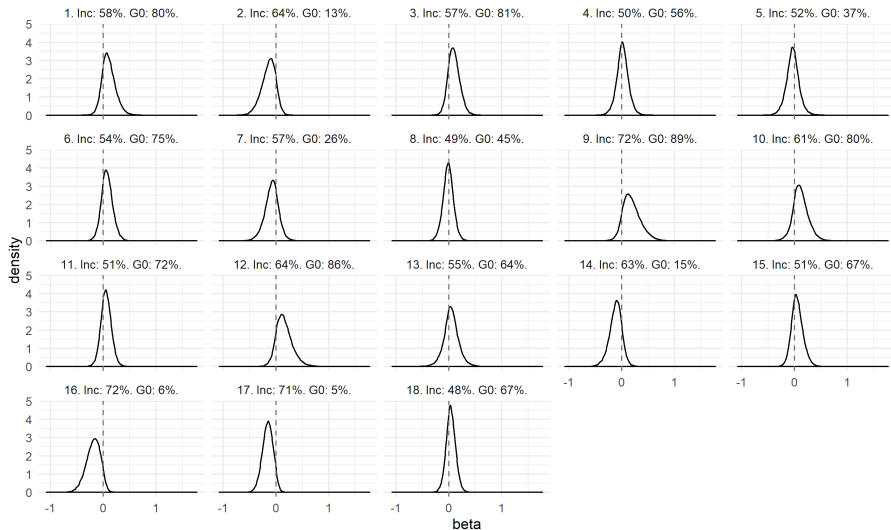
Micro 2: 20.0% and 18.7%

Macro 2 56.1% and 55.2% (57th and 58th most included out of 62)

Finance 1: 71.0% and 48.4% (3rd and 18th most included out of 18)

I. Adding meaningless variables

Finance 1 data set ($n = 68$): Here x_{17} and x_{18} are meaningless.



Similar shapes: β_{18} and $(\beta_4, \beta_5, \beta_{15})$.

High inclusion: x_{17} included 71% of times.

II. Fatter tails via Student's t

New prior:

$$\beta_i | \sigma^2, \gamma^2, \lambda_i^2, q \sim \begin{cases} N(0, \sigma^2 \gamma^2 \lambda_i^2) & \text{with prob. } q \\ 0 & \text{with prob. } 1 - q \end{cases}$$
$$\lambda_i^2 \sim IG\left(\frac{\nu}{2}, \frac{\nu}{2}\right).$$

Therefore, β_i follows a Student's t distribution:

$$\beta_i | \sigma^2, \gamma^2, q \sim \begin{cases} t_\nu(0, \sigma^2 \gamma^2) & \text{with prob. } q \\ 0 & \text{with prob. } 1 - q \end{cases}$$

where

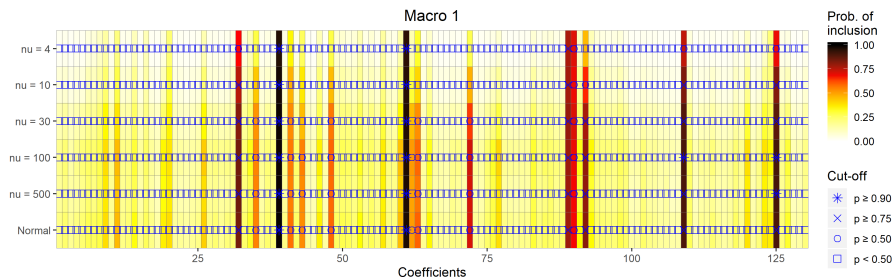
$$V(\beta_i | \sigma^2, \gamma^2, q) = \frac{\nu}{\nu - 2} \sigma^2 \gamma^2$$

II. Fatter tails via Student's t - Macro 1

x_{72} and x_{90} are both relevant for $\nu > 10$.

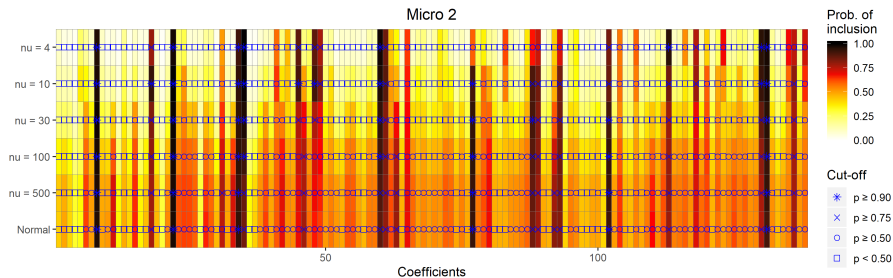
Only x_{90} for $\nu \leq 10$ (**sparsity reemerges**).

Prob. inclusion \downarrow as $\nu \uparrow$.



Argument: Spike-and-Slab, as originally defined, induces selection and shrinkage, since for $\nu = 4$ only 7 of 130 available predictors are relevant - that is, included more than 50% of the times.

II. Fatter tails via Student's t - Micro 2



Gaussian: no pattern of variable selection.

106 of 138 predictors are selected more than 50% of the times.

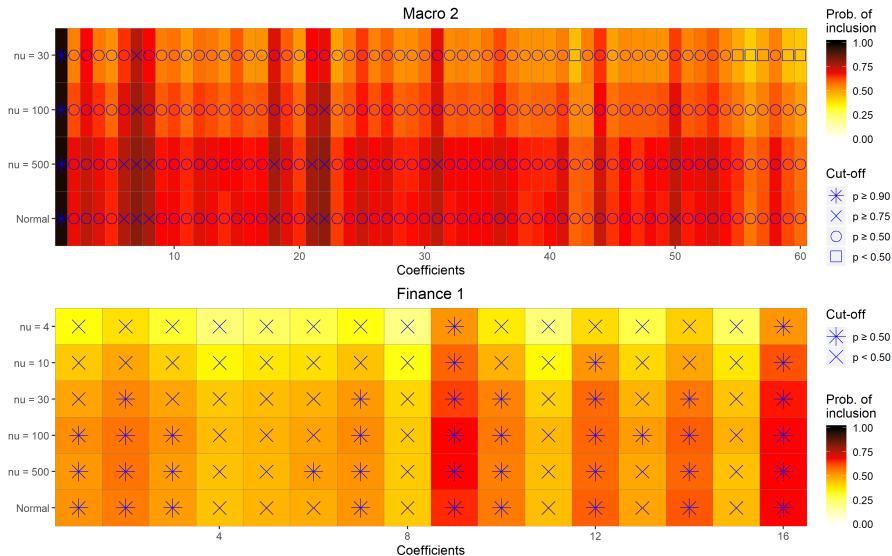
Student's t : Sparsity in action.

For $\nu = 4$, only 30 predictors are selected.

For $\nu = 10$, only 34 predictors are selected.

II. Fatter tails via Student's t - Macro 2 & Finance 1

Similarity across ν



III. More simulators

For observation $i = 1, \dots, n = 68$ and predictor $j = 1, \dots, k = 16$, we simulate

$$x_{ij} \sim N(0, 1) \quad \text{and} \quad \varepsilon_i^* \sim N(0, 1)$$

We also fix $\beta_1 = -0.86$, $\beta_2 = 0.64$ and $\beta_3 = 0.89$, while the response variable is:

$$y_i^{(s)} = \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \sigma_\varepsilon^{(s)} \varepsilon_i^*,$$

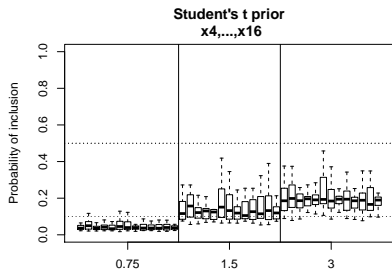
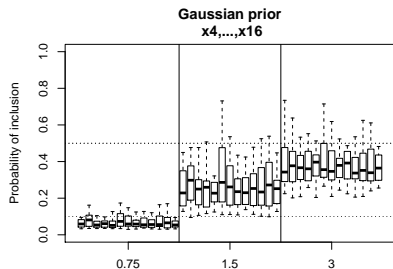
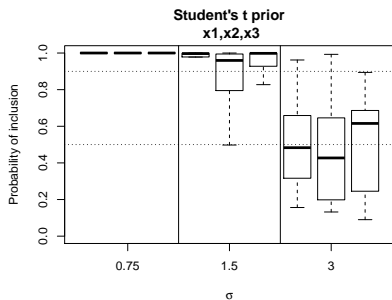
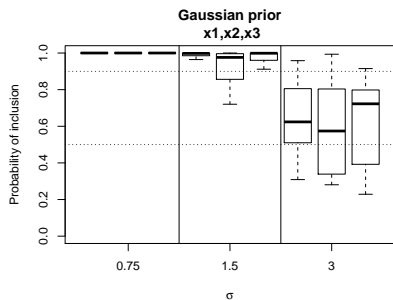
and $\sigma_\varepsilon^{(s)} = 0.75s$, for $s = 1, 2, 3$.

The prior for β are Gaussian or Student's t with $\nu = 4$ degrees of freedom.

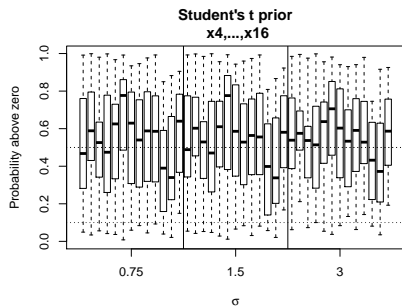
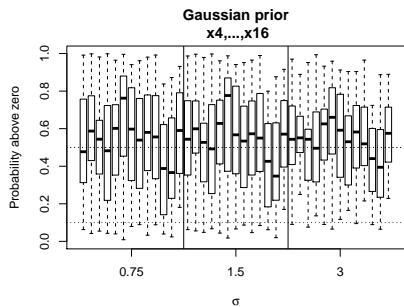
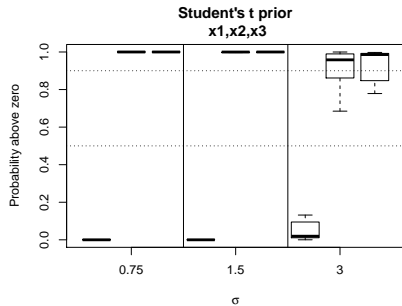
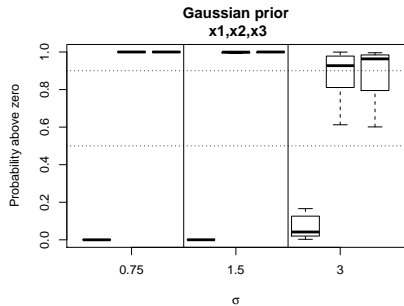
We replicate the above simulation $R = 20$ times.

III. Probability of inclusion

- ▶ $\sigma \uparrow$: inclusion of x_1, x_2, x_3 decreases. More so for the Student's t case.
- ▶ $\sigma \uparrow$: inclusion of x_4, \dots, x_{16} increases. More so for the Gaussian case.

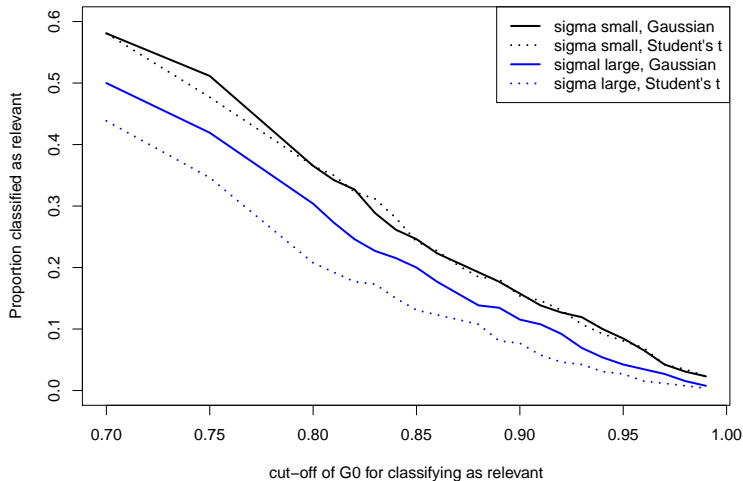


III. Probability above zero



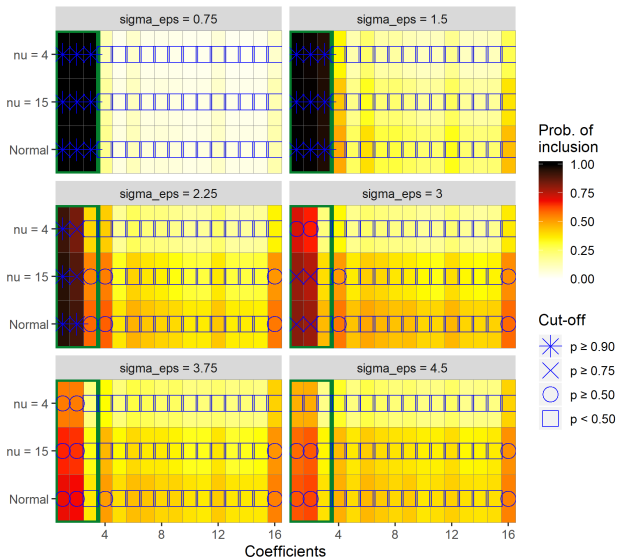
III. Proportion of $\beta_4, \dots, \beta_{16}$ classified as relevant

For σ large, Student's t prior performs better at shrinking towards zero.



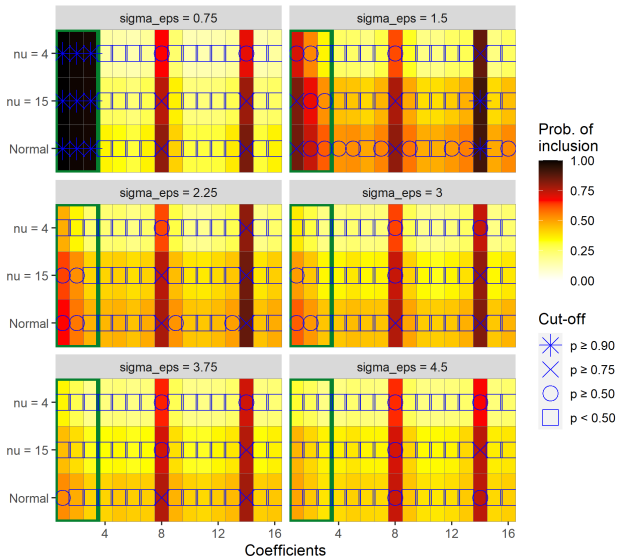
Independent predictors

$\beta_1 = -0.86, \beta_2 = 0.64, \beta_3 = 0.89$



Correlated predictors $\rho = 0.8$

$\beta_1 = -0.86, \beta_2 = 0.64, \beta_3 = 0.89$



Outline

Sparsity in macro, micro and finance

GLP approach

Our contribution

Brief review of sparsity in linear models

- Ridge and lasso regressions

- Spike and slab model (or SMN model)

- SSVS and scaled SSVS priors

- R package `Bayeslm`

Recalling GLP main remarks

- Inclusion and tail probabilities

Our experiments

- I. Adding meaningless variables

- II. Fatter tails via Student's t

- III. More simulations

Final remarks

Final remarks

- ▶ Their illusion resides, in our view, mainly on the anticipated expectation that sparsity would be as obvious in Economic applications as it has been in various other fields, such as genomics and machine-learning (ML) applications.
- ▶ A natural extension (they touch it briefly in the paper) is to allow for common factors (or principal components) along with the multiple regressors. See, for instance, Hahn, He and Lopes (2018) **Bayesian factor model shrinkage for linear IV regression with many instruments**, *Journal of Business and Economic Statistics*, 2018, 36(2), 278-287.
- ▶ The data sets used in GLP are diverse, but still represent a modest increase compared to the vertiginous sizes of data sets found elsewhere, say in the ML literature.
- ▶ Financial econometrics and micro-econometrics might be riper now for implementations of such sparsity and shrinkage inducing priors.

References

- Sylvia Frühwirth-Schnatter and Hedibert F. Lopes. Sparse Bayesian factor analysis when the number of factors is unknown. Technical report, 2018.
- Edward I George and Robert E McCulloch. Variable selection via gibbs sampling. *Journal of the American Statistical Association*, 88(423):881–889, 1993.
- Jim Griffin and Philip Brown. Inference with normal-gamma prior distributions in regression problems. *Bayesian Analysis*, 5(1):171–188, 2010.
- Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Hemant Ishwaran and J Sunil Rao. Spike and slab variable selection: frequentist and bayesian strategies. *Annals of Statistics*, pages 730–773, 2005.
- Gregor Kastner, Sylvia Frühwirth-Schnatter, and Hedibert F. Lopes. Efficient Bayesian inference for multivariate factor stochastic volatility models. *Journal of Computational and Graphical Statistics*, 26: 905–917, 2017.
- F. B. Lempers. *Posterior Probabilities of Alternative Linear Models*. Rotterdam University Press, 1988.
- T. J. Mitchell and J. J. Beauchamp. Bayesian variable selection in linear regression (with discussion). *Journal of the American Statistical Association*, 83:1023–1036, 1988.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.