

Bayesian Variable Selection - Computational methods

Jim Griffin

University College London



Plan

- Basic search and MCMC methods
- High-dimensional / large p problems
- R packages

MCMC for variable selection

The model space $\Gamma = \{0, 1\}^P$ is a discrete space (e.g. there are 2^p possible regression models with p variables)

As p becomes larger, the number of models grows exponentially and it is difficult and time-consuming to calculate the posterior model probabilities for all models if $p > 30$

Approaches

There are two main approaches to this problem

- Approximate the posterior distribution for a subset of the $\Gamma' \subseteq \Gamma$ by

$$p(\gamma|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \gamma) p(\gamma)}{\sum_{\gamma' \in \Gamma'} p(\mathbf{y}|\mathbf{X}, \gamma') p(\gamma')}, \quad \gamma \in \Gamma'$$

if $p(\mathbf{y}|\mathbf{X}, \gamma)$ is available analytically

- An alternative is to sample models $\gamma^{(1)}, \dots, \gamma^{(N)}$ according to the posterior model probabilities $p(\gamma|\mathbf{X}, \mathbf{y})$ and approximate

$$p(\gamma|\mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N \mathbf{I}(\gamma^{(i)} = \gamma), \quad \gamma \in \Gamma$$

Stochastic Search Variable Selection (SSVS)

In stochastic search variable selection, a heuristic algorithm is developed to find “most” of the model with highest posterior probabilities to use as Γ' .

A prominent example is Shotgun Stochastic Search (SSS) (Hans et al., 2007) where at each iteration a **neighbourhood** is constructed around the current model by either: including an extra variable in the model, removing a variable from the model, swapping a variable in the model with one outside.

Sample a new model in proportion to its (unnormalized) posterior model probability

Comments

- The set Γ' can include all models visited in a neighbourhood
- Calculations across neighbourhoods can be parallelized
- Bayesian Adaptive Sampling (Clyde et al., 2011)
samples without replacement from the posterior using a representation of the posterior distribution as a binary tree
- Scaling to large p leads to large neighbourhoods. Shin et al. (2018) and Li et al. (2023) discuss computational methods and screening methods to control computational cost

Gibbs sampler

A Gibbs sampler can be run on γ with full conditionals

$$\begin{aligned} \pi(\gamma_k = 1 | \gamma_{-k}) \\ = \frac{p(\mathbf{y} | \mathbf{X}, \gamma_k = 1, \gamma_{-k}) p(\gamma_k = 1, \gamma_{-k})}{p(\mathbf{y} | \mathbf{X}, \gamma_k = 1, \gamma_{-k}) p(\gamma_k = 1, \gamma_{-k}) + p(\mathbf{y} | \mathbf{X}, \gamma_k = 0, \gamma_{-k}) p(\gamma_k = 1, \gamma_{-k})} \end{aligned}$$

These one-at-a-time updates can lead to slow mixing

Metropolis-Hastings sampler

The standard approach is Metropolis-Hastings sampling with the following proposal distribution for the probability of proposing γ' if you are currently at γ . Let p_γ be the number of included variables in model γ then, for $a < 0.5$,

- Add – with probability a , propose to **add** a variable uniformly at random from those **excluded** in model γ .
- Delete – with probability a , propose to **delete** a variable uniformly at random from those **included** in model γ .
- Swap – with probability $1 - 2a$, propose to add a variable uniformly at random from those excluded in model γ and propose to delete a variable uniformly at random from those included in model γ .

Acceptance probability

The acceptance probability is

$$\alpha(\gamma, \gamma') = \min \left\{ 1, \frac{p(\mathbf{y}|\mathbf{X}, \gamma')p(\gamma')q(\gamma', \gamma)}{p(\mathbf{y}|\mathbf{X}, \gamma)p(\gamma)q(\gamma, \gamma')} \right\}$$

For add, $q(\gamma, \gamma') = a \frac{1}{p - p_\gamma}$

For delete, $q(\gamma, \gamma') = a \frac{1}{p_\gamma}$

For swap, $q(\gamma, \gamma') = (1 - 2a) \frac{1}{(p - p_\gamma)p_\gamma}$

Large p models

The availability of data sets with p in the thousands has led to a focus on performance of MCMC algorithms.

Yang et al. (2016) show that an ADS sampler is rapidly mixing (*i.e.* the mixing time is approximately $p n s_0^2 \log p$)

This has motivated new samplers which make use the following ideas to improve proposals and mixing of the chains.

- **Informed proposals**: can we improve proposals by including likelihood information?
- **Adaptation**: can we use previously sampled values to improve proposals?

Locally informed proposals

Zanella (2020) introduced locally informed proposals. We use the following neighbourhood-based definition

$$q_{N(\gamma),g}(\gamma') = \frac{g\left(\frac{p(\gamma'|\text{Data})}{p(\gamma|\text{Data})}\right)}{Z_g(\gamma)}, \quad \gamma' \in N(\gamma)$$

where

- $N(\gamma)$ is a neighbourhood of γ (e.g. 1 Hamming distance ball)
- g is a **balancing function** for which $g(t) = t g(1/t)$
e.g. $g(t) = \min\{1, t\}$ or $g(t) = t/(1+t)$

■

$$Z_g(\gamma) = \sum_{\tilde{\gamma} \in N(\gamma)} g\left(\frac{p(\tilde{\gamma}|\text{Data})}{p(\gamma|\text{Data})}\right)$$

Zanella (2020) argues that

- locally balanced proposals are optimal asymptotically in the number of regressors
- the improvement over a locally unbalanced function increases as the posterior distribution becomes rougher

Zanella (2020) also derives the optimal $g = \frac{t}{1+t}$ for a Add-Delete-Swap neighbourhood and an independent posterior

Weighted Tempered Gibbs Sampling (Zanella and Roberts, 2019)

Zanella and Roberts (2019) introduce a general theory for sampling from high-dimensional posteriors using Gibbs sampling with tempered posteriors and importance sampling

WTGS algorithm for BVS: at each iteration of the Markov chain do

- (a) sample i from $\{1, \dots, p\}$ proportionally to

$$p_i(\gamma) = \frac{p(\gamma_i = 1 | \gamma_{-i}, Y) + k/p}{2p(\gamma_i | \gamma_{-i}, Y)};$$

- (b) switch γ_i to $1 - \gamma_i$;
(c) weight the new state γ with a weight $Z(\gamma)^{-1}$ where $Z(\gamma) \propto \sum_{i=1}^p p_i(\gamma)$

Adaptively Scaled Independence (ASI) proposal (Griffin et al., 2021)

Let π_1, \dots, π_p be the PIPs and propose $\gamma' = (\gamma'_1, \dots, \gamma'_p)$ independently using

	$\gamma'_k = 0$	$\gamma'_k = 1$
$\gamma_k = 0$	$1 - \omega A_k$	ωA_k
$\gamma_k = 1$	$1 - \omega D_k$	ωD_k

where

- $A_k = \min \{1, \pi_k / (1 - \pi_k)\}$ and $D_k = \min \{1, (1 - \pi_k) / \pi_k\}$
- ω allows us to control the number of variables changed

- If γ_k 's are independent under the posterior then this is optimal proposal with $\omega = 1$
- The PIPs can be approximated using the previous MCMC output
- ω can be chosen to control the overall acceptance rate

Adaptive Random Neighbourhood Informed (ARNI) (Liang et al., 2022)

To scale informed proposals to large p , we use a **random neighbourhood** which allows us to

- control the distribution of the size of $N(\gamma)$
- favour variables which improve the quality of the neighbourhood

Algorithm

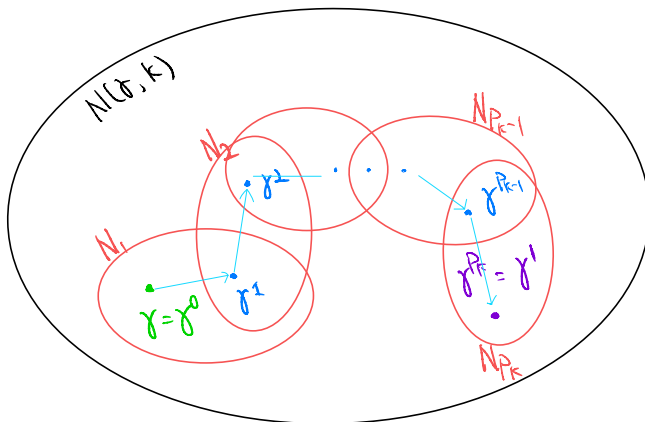
- Neighbourhood generation: Draw $N \sim p(N|\gamma)$.
- (Informed) within neighbourhood proposal: $\gamma' \sim q_{N,g}$ where

$$q_{N,g}(\gamma') = \frac{1}{Z_{N,g}(\gamma)} g \left(\frac{p(\gamma' | \text{Data}) p(N|\gamma') q_N(\gamma)}{p(\gamma | \text{Data}) p(N|\gamma) q_N(\gamma')} \right) q_N(\gamma')$$

- Accept/reject step

$$\alpha = \min \left\{ 1, \frac{p(\gamma' | \text{Data}) p(N|\gamma') q_N(\gamma)}{p(\gamma | \text{Data}) p(N|\gamma) q_N(\gamma')} \right\} = \min \left\{ 1, \frac{Z_{N,g}(\gamma')}{Z_{N,g}(\gamma)} \right\}$$

PARNI sampler (Liang et al., 2022)



Neighbourhood generation (Liang et al., 2022)

Choose a set κ of variables to change in the model using the ASI proposal

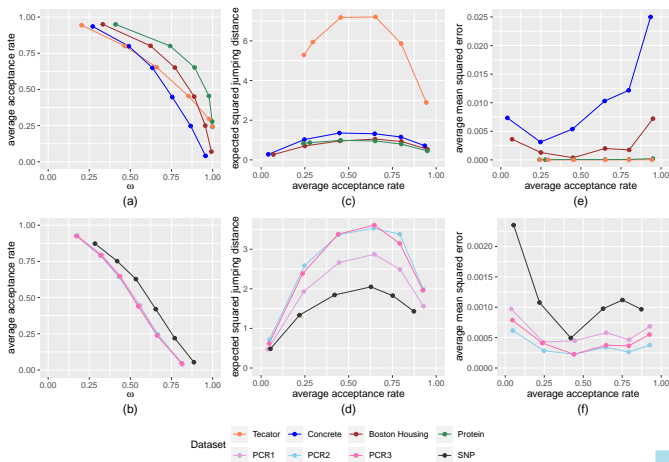
	$\kappa_k = 0$	$\kappa_k = 1$
$\gamma_k = 0$	$1 - \omega A_k$	ωA_k
$\gamma_k = 1$	$1 - \omega D_k$	ωD_k

ω is tuned using the Robbins-Monro or Kiefer-Wolfowitz algorithms with target acceptance rate 0.65, which is similar to informed proposals on continuous spaces:

- Metropolis-adjusted Langevin algorithm (MALA) – 0.57 (Roberts and Rosenthal, 1998)
- Hamiltonian Monte Carlo – 0.65 (Beskos et al, 2013)

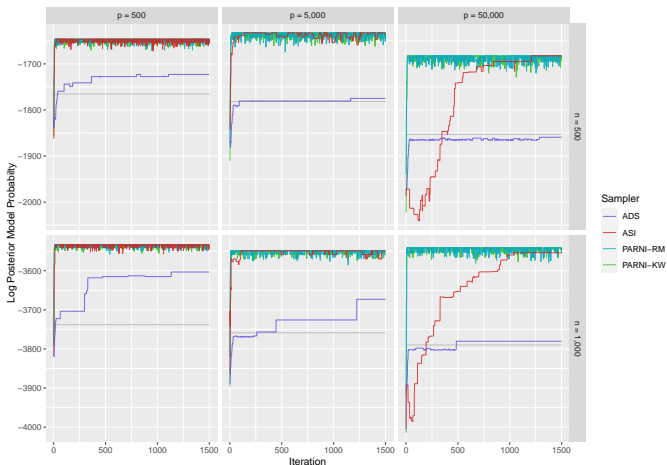
Optimal value of ω

Liang, Livingstone, Griffin (Stats and Computing, 2022)



Convergence – Simulated data

Liang, Livingstone, Griffin (Stats and Computing, 2022)



Relative average mean squared errors for inclusion variables (difference to Add-Delete-Swap)

dataset	n	p	ASI	HBS	WTGS	LIT	PARNI KW	PARNI RM
Tecator	172	100	-0.50	0.31	-0.25	0.55	0.09	0.09
Boston	506	104	0.73	-1.14	-0.72	-0.57	-0.91	-0.90
Housing								
Concrete	1030	79	0.09	0.45	-0.17	-0.15	0.45	1.00
Protein	96	88	-0.47	1.54	-1.05	-0.31	-0.47	-0.64
PCR1	60	22,575	-0.91	0.52	-1.15	-0.20	-1.17	-0.98
PCR2	60	22,575	-0.30	1.12	-0.56	0.01	-0.79	-0.87
PCR3	60	22,575	-0.11	1.37	-0.42	0.19	-0.59	-0.37
SNP	993	79,748	-0.57	0.39	-0.83	-0.97	-1.45	-1.31

R package

The BAS package (Clyde et al., 2011) allows inference in a range of models including linear regression, models with interactions, and generalized linear models using the standard `lm`-type model definition

It includes a wide range of priors, computational methods, graphical and numerical displays

The BayesVarSel package (Garcia-Donato and Forte, 2018) is a package for Gibbs sampling with a linear regression and a wide range of priors

- Clyde, M. A., Ghosh, J. and Littman, M. L. (2011). Bayesian adaptive sampling for variable selection and model averaging, *JCGS* **20**: 80–101.
- Garcia-Donato, G. and Forte, A. (2018). Bayesian testing, variable selection and model averaging in linear models using R with BayesVarSel, *The R Journal* **10**(1): 329.
URL: <https://journal.r-project.org/archive/2018/RJ-2018-021/index.html>
- Griffin, J. E., Latuszynski, K. G. and Steel, M. F. J. (2021). In search of lost mixing time: adaptive Markov chain Monte Carlo schemes for Bayesian variable selection with very large p , *Biometrika* **108**: 53–69.
- Hans, C., Dobra, A. and West, M. (2007). Shotgun stochastic search for “large p ” regression, *JASA* **102**: 507–516.
- Li, D., Dutta, S. and Roy, V. (2023). Model based screening embedded Bayesian variable selection for ultra-high dimensional settings, *JCGS* **32**: 61 – 73.
- Liang, X., Livingstone, S. and Griffin, J. E. (2022). Adaptive random neighbourhood informed Markov chain Monte Carlo for high-dimensional Bayesian variable selection, *Stat. and Comp.* **32**: 84.
- Shin, M., Bhattacharya, A. and Johnson, V. E. (2018). Scalable Bayesian variable selection using nonlocal prior densities in ultrahigh-dimensional settings, *Statistica Sinica* **28**: 1053–1078.

- Yang, Y., Wainwright, M. J. and Jordan, M. I. (2016). On the computational complexity of high-dimensional Bayesian variable selection, *The Annals of Statistics* **44**: 2497–2532.
- Zanella, G. (2020). Informed proposals for local MCMC in discrete spaces, *JASA* **115**: 852–865.
- Zanella, G. and Roberts, G. (2019). Scalable importance tempering and Bayesian variable selection, *JRSS B* **81**: 489–517.