



This is a repository copy of *On spike-and-slab priors for Bayesian equation discovery of nonlinear dynamical systems via sparse linear regression*.

White Rose Research Online URL for this paper:
<https://eprints.whiterose.ac.uk/173869/>

Version: Supplemental Material

Article:

Nayek, R. orcid.org/0000-0003-4277-8382, Fuentes, R., Worden, K. et al. (1 more author) (2021) On spike-and-slab priors for Bayesian equation discovery of nonlinear dynamical systems via sparse linear regression. *Mechanical Systems and Signal Processing*, 161. 107986. ISSN 0888-3270

<https://doi.org/10.1016/j.ymssp.2021.107986>

Article available under the terms of the CC-BY-NC-ND licence
(<https://creativecommons.org/licenses/by-nc-nd/4.0/>).

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Supplementary Material

On spike-and-slab priors for Bayesian equation discovery of nonlinear dynamical systems via sparse linear regression

R. Nayek, R. Fuentes, K. Worden, E.J. Cross

This document comprises the derivation of the conditional probability distributions in the Gibbs sampling scheme used for Discontinuous Spike-and-Slab (DSS) priors. The spike-and-slab prior has a hierarchical form that can be expressed in the form of a directed acyclic graph (DAG), and hence, it will be useful to understand the play of conditional independencies between the random variables in a DAG, before looking to derive the conditional probability distributions of the Gibbs sampler.

1. Conditional independence in DAGs

Given three random variables, say X , Y and Z ,

$$X \perp Y \mid Z \quad (1)$$

implies that X and Y are probabilistically *independent* given $Z = z$. Conditional independencies between random variables are easily visualised using DAGs, and whether or not two random variables (depicted by the nodes in a DAG) are conditionally independent is decided based on their structure of their connection i.e. the directed edges of the DAG. Figure 1 shows three different types of DAG structures.

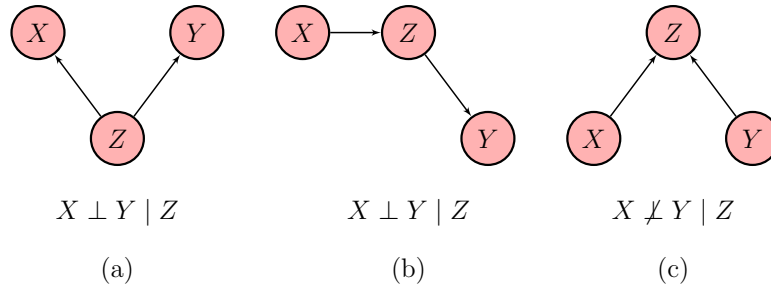


Figure 1: Three different types of directed acyclic graph structures.

In the first two DAG structures, 1(a) and 1(b), the conditioning on Z renders Y independent of X , and therefore, graph structures similar to cases (a) and (b) would imply $X \perp Y \mid Z$. For DAG structure 1(c), X and Y are marginally independent of each other, that is $X \perp Y$, however, when conditioned on Z , X and Y become dependent, i.e. $X \not\perp Y \mid Z$. The notion of conditional independence in DAGs would greatly aid in finding the conditional probability distributions feeding the Gibbs sampler.

2. Derivation of the conditional probability distributions of Gibbs sampler for DSS priors

2.1. DSS prior model specification

The DSS prior model used for linear regression with dictionary $\mathbf{D} \in \mathbb{R}^{N \times P}$ is summarised by the following set of relations:

$$p(\mathbf{y} | \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(\mathbf{D}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_N) \quad (2a)$$

$$p(\boldsymbol{\theta} | \mathbf{z}) = p_{\text{slab}}(\boldsymbol{\theta}_r) \prod_{i: z_i=0} p_{\text{spike}}(\theta_i) \quad (2b)$$

$$p_{\text{spike}}(\theta_i) = \delta_0 \quad (2c)$$

$$p_{\text{slab}}(\boldsymbol{\theta}_r | \sigma^2, v_s) = \mathcal{N}(\mathbf{0}, \sigma^2 v_s \mathbf{A}_{0,r}) \quad (2d)$$

$$p(v_s) = \mathcal{IG}(a_v, b_v) \quad (2e)$$

$$p(z_i | p_0) = \text{Bern}(p_0) \quad (2f)$$

$$p(p_0) = \text{Beta}(a_p, b_p) \quad (2g)$$

$$p(\sigma^2) = \mathcal{IG}(a_\sigma, b_\sigma) \quad (2h)$$

Here, $\boldsymbol{\theta}_r \in \mathbb{R}^{r \times 1}$ denotes the set of components of $\boldsymbol{\theta}$ for which the corresponding z_i s take values 1. The relations between the random variables can be visualised as in Figure 2.

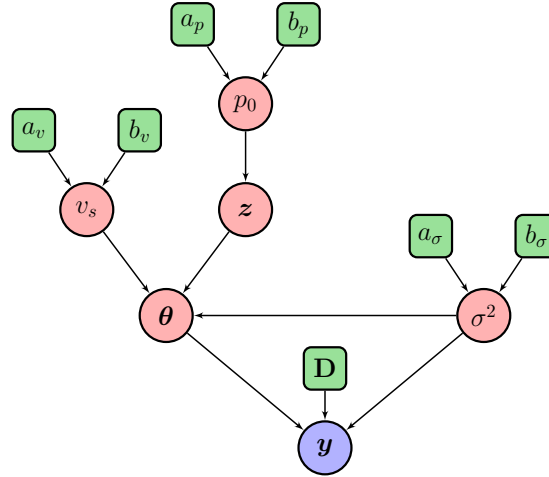


Figure 2: Graph of the hierarchical spike-and-slab model for linear regression; the variables in circles represent random variables, while those in squares represent deterministic parameters.

Due to the DAG structure, the joint probability distribution over all random variables factorises as,

$$p(\mathbf{y}, \boldsymbol{\theta}, \mathbf{z}, p_0, v_s, \sigma^2) = p(\mathbf{y} | \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} | \mathbf{z}, v_s, \sigma^2) p(\mathbf{z} | p_0) p(p_0) p(v_s) p(\sigma^2) \quad (3)$$

Note the conditioning on deterministic hyperparameters $a_v, b_v, a_p, b_p, a_\sigma, b_\sigma$ and the dictionary \mathbf{D} has been suppressed. The Gibbs sampler needs closed-form expressions of the conditional probability distributions of each random variable given the data and all other random variables. However, due to the DAG structure, the conditional probability distributions of the random variables become conditionally independent of certain variables:

$$p(p_0 | \mathbf{y}, \boldsymbol{\theta}, \mathbf{z}, v_s, \sigma^2) = p(p_0 | \mathbf{z}) \quad (4a)$$

$$p(v_s | \mathbf{y}, \boldsymbol{\theta}, \mathbf{z}, p_0, \sigma^2) = p(v_s | \boldsymbol{\theta}, \mathbf{z}, \sigma^2) \quad (4b)$$

$$p(\sigma^2 | \mathbf{y}, \boldsymbol{\theta}, \mathbf{z}, p_0, v_s) = p(\sigma^2 | \mathbf{y}, \boldsymbol{\theta}, \mathbf{z}, v_s) \quad (4c)$$

$$p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{z}, p_0, v_s, \sigma^2) = p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{z}, v_s, \sigma^2) \quad (4d)$$

$$p(\mathbf{z} | \mathbf{y}, \boldsymbol{\theta}, p_0, v_s, \sigma^2) = p(\mathbf{z} | \boldsymbol{\theta}, p_0, v_s, \sigma^2) \quad (4e)$$

At this point, it must be mentioned that the Dirac-delta spike distribution in DSS priors creates a nuisance in the Gibbs sampling scheme; it prevents the creation of an irreducible Markov chain required for the chain to reach a stationary distribution. Specifically, when $z_i = 0$, the new value for θ_i will be $\theta_i = 0$, which will

in turn result in $z_i = 0$ in a new draw, and the process keeps repeating itself. In other words, the chain has absorbing states. This problem can be avoided by *marginalising* or integrating over $\boldsymbol{\theta}$ to get rid of the Dirac-delta function. Therefore, in the derivations of the conditional probability distributions for Eqs. 4c and 4e, the procedure of marginalising over $\boldsymbol{\theta}$ is performed to obtain an irreducible chain. As a result of marginalisation over $\boldsymbol{\theta}$, the conditional probabilities Eqs. 4c and 4e become dependent upon \mathbf{y} , and their relations can be re-expressed as,

$$p(\sigma^2 \mid \mathbf{y}, \boldsymbol{\theta}, \mathbf{z}, p_0, v_s) = p(\sigma^2 \mid \mathbf{y}, \mathbf{z}, v_s) \quad (5a)$$

$$p(\mathbf{z} \mid \mathbf{y}, \boldsymbol{\theta}, p_0, v_s, \sigma^2) = p(\mathbf{z} \mid \mathbf{y}, p_0, v_s, \sigma^2) \quad (5b)$$

2.2. Derivations of conditional probability distributions

In this section, the conditional sampling distributions used in Gibbs sampling with DSS priors are derived.

2.2.1. Sampling distribution of p_0

$$\begin{aligned} p(p_0 \mid \mathbf{z}) &\propto p(\mathbf{z} \mid p_0) p(p_0) \\ &\propto \left(\prod_{i=1}^P p_0^{z_i} (1-p_0)^{1-z_i} \right) p_0^{a_p-1} (1-p_0)^{b_p-1} \quad [\because z_i \text{'s are independent}] \\ &\propto p_0^{a_p+s_z} (1-p_0)^{b_p+P-s_z} \quad \left[\text{where } s_z = \sum_{i=1}^P z_i \right] \\ &\boxed{p_0 \mid \mathbf{z} \sim \text{Beta}(a_p + s_z, b_p + P - s_z)} \end{aligned} \quad (6)$$

2.2.2. Sampling distribution of v_s

$$\begin{aligned} p(v_s \mid \boldsymbol{\theta}, \mathbf{z}, \sigma^2) &\propto p(\boldsymbol{\theta} \mid \mathbf{z}, v_s, \sigma^2) p(v_s) p(\mathbf{z}) p(\sigma^2) \\ &\propto p(\boldsymbol{\theta} \mid \mathbf{z}, v_s, \sigma^2) p(v_s) \quad [\because p(\mathbf{z}) \text{ and } p(\sigma^2) \text{ are constants w.r.t. } v_s] \\ &\propto \mathcal{N}(\boldsymbol{\theta}_r \mid \mathbf{0}, \sigma^2 v_s \mathbf{A}_{0,r}) \mathcal{IG}(a_v, b_v) \\ &\propto \frac{1}{(v_s \sigma^2)^{r/2} |\mathbf{A}_{0,r}|^{1/2}} \exp\left(-\frac{\boldsymbol{\theta}_r^T \mathbf{A}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2 v_s}\right) v_s^{-a_v-1} \exp\left(-\frac{b_v}{v_s}\right) \\ &\propto v_s^{-(a_v + \frac{r}{2})-1} \exp\left(-\frac{b_v + \frac{\boldsymbol{\theta}_r^T \mathbf{A}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2}}{v_s}\right) \\ &\boxed{v_s \mid \boldsymbol{\theta}, \mathbf{z}, \sigma^2 \sim \mathcal{IG}\left(a_v + \frac{r}{2}, b_v + \frac{\boldsymbol{\theta}_r^T \mathbf{A}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2}\right)} \end{aligned} \quad (7)$$

Note, r is the number of components of $\boldsymbol{\theta}$ that fall in the slab and is equal to the number of non-zero components of \mathbf{z} . Since z_i s are binary variables, $r = s_z = \sum_{i=1}^P z_i$. When all z_i s are zero, v_s is sampled from the prior $\mathcal{IG}(a_v, b_v)$.

2.2.3. Sampling distribution of σ^2

The derivation of conditional distribution of σ^2 will involve integrating out the parameter $\boldsymbol{\theta}$ to avoid dealing with the Dirac-delta spike distribution (as mentioned before).

$$\begin{aligned} p(\sigma^2 \mid \mathbf{y}, \mathbf{z}, v_s) &\propto \int p(\mathbf{y}, \boldsymbol{\theta}, \mathbf{z}, v_s, \sigma^2) d\boldsymbol{\theta} \\ &\propto \left(\int p(\mathbf{y} \mid \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} \mid \mathbf{z}, v_s, \sigma^2) d\boldsymbol{\theta} \right) p(\mathbf{z}) p(v_s) p(\sigma^2) \\ &\propto \left(\int p(\mathbf{y} \mid \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} \mid \mathbf{z}, v_s, \sigma^2) d\boldsymbol{\theta} \right) p(\sigma^2) \quad [\because p(\mathbf{z}) \text{ and } p(v_s) \text{ are constants w.r.t. } \sigma^2] \end{aligned}$$

Upon integration, the components of $\boldsymbol{\theta}$ belonging to the spike distribution are evaluated at zero and the rest of the components of $\boldsymbol{\theta}$ that belong to the slab (i.e. $\boldsymbol{\theta}_r$) remains to be integrated. Expanding the integrand $p(\mathbf{y} \mid \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} \mid \mathbf{z}, v_s, \sigma^2)$, with $\boldsymbol{\theta}$ now replaced by $\boldsymbol{\theta}_r$, one obtains,

$$\frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{y} - \mathbf{D}_r \boldsymbol{\theta}_r)^T (\mathbf{y} - \mathbf{D}_r \boldsymbol{\theta}_r)}{2\sigma^2}\right) \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2}}{(2\pi v_s \sigma^2)^{r/2}} \exp\left(-\frac{\boldsymbol{\theta}_r^T \mathbf{A}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2 v_s}\right)$$

Simplifying and rearranging the terms involving $\boldsymbol{\theta}_r$ leads to,

$$\begin{aligned} &\Rightarrow \frac{1}{(2\pi\sigma^2)^{N/2}} \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2}}{(2\pi v_s \sigma^2)^{r/2}} \exp\left(-\frac{\mathbf{y}^T \mathbf{y} + \boldsymbol{\theta}_r^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta}_r - 2\boldsymbol{\theta}_r^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}}{2\sigma^2}\right) \\ &\Rightarrow \frac{1}{(2\pi\sigma^2)^{N/2}} \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2}}{(2\pi v_s \sigma^2)^{r/2}} \exp\left(-\frac{(\boldsymbol{\theta}_r - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta}_r - \boldsymbol{\mu})}{2\sigma^2}\right) \exp\left(-\frac{(\mathbf{y}^T \mathbf{y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{2\sigma^2}\right) \end{aligned}$$

where $\boldsymbol{\Sigma}^{-1} = (\mathbf{D}_r^T \mathbf{D}_r + v_s^{-1} \mathbf{A}_{0,r}^{-1})$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma} \mathbf{D}_r^T \mathbf{y}$. On integrating out $\boldsymbol{\theta}_r$ from the above expression, the conditional distribution of σ^2 reduces to:

$$\begin{aligned} p(\sigma^2 \mid \mathbf{y}, \mathbf{z}, v_s) &\propto \frac{1}{(\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{y}^T \mathbf{y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{2\sigma^2}\right) p(\sigma^2) \\ &\propto \frac{1}{(\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{y}^T \mathbf{y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{2\sigma^2}\right) (\sigma^2)^{-a_\sigma - 1} \exp\left(-\frac{b_\sigma}{\sigma^2}\right) \\ &\propto (\sigma^2)^{-(a_\sigma + \frac{N}{2}) - 1} \exp\left(-\frac{b_\sigma + \frac{1}{2} (\mathbf{y}^T \mathbf{y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{\sigma^2}\right) \end{aligned}$$

$$\boxed{\sigma^2 \mid \mathbf{y}, \mathbf{z}, v_s \sim \mathcal{IG}\left(a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\mathbf{y}^T \mathbf{y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})\right)} \quad (8)$$

Note that when all z_i s are equal to zero, σ^2 is sampled from $\mathcal{IG}(a_\sigma + \frac{N}{2}, b_\sigma + \frac{1}{2} (\mathbf{y}^T \mathbf{y}))$.

2.2.4. Sampling distribution of $\boldsymbol{\theta}$

The components of $\boldsymbol{\theta}$ that correspond to $z_i = 0$ are set to zero (as they belong to the spike Dirac-delta distribution). The conditional distribution for the rest of the components belonging to the slab, represented by $\boldsymbol{\theta}_r$, can be derived as follows:

$$\begin{aligned} p(\boldsymbol{\theta}_r \mid \mathbf{y}, v_s, \sigma^2) &\propto \mathcal{N}(\mathbf{y} \mid \mathbf{D}_r \boldsymbol{\theta}_r, \sigma^2 \mathbf{I}_N) \mathcal{N}(\boldsymbol{\theta}_r \mid \mathbf{0}, \sigma^2 v_s \mathbf{A}_{0,r}) \\ &\propto \exp\left(-\frac{(\mathbf{y} - \mathbf{D}_r \boldsymbol{\theta}_r)^T (\mathbf{y} - \mathbf{D}_r \boldsymbol{\theta}_r)}{2\sigma^2}\right) \exp\left(-\frac{\boldsymbol{\theta}_r^T \mathbf{A}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2 v_s}\right) \\ &\propto \exp\left(-\frac{\mathbf{y}^T \mathbf{y} + \boldsymbol{\theta}_r^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta}_r - 2\boldsymbol{\theta}_r^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}}{2\sigma^2}\right) \\ &\propto \exp\left(-\frac{(\boldsymbol{\theta}_r - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta}_r - \boldsymbol{\mu})}{2\sigma^2}\right) \end{aligned}$$

where $\boldsymbol{\Sigma}^{-1} = (\mathbf{D}_r^T \mathbf{D}_r + v_s^{-1} \mathbf{A}_{0,r}^{-1})$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma} \mathbf{D}_r^T \mathbf{y}$.

$$\boxed{\boldsymbol{\theta}_r \mid \mathbf{y}, v_s, \sigma^2 \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \boldsymbol{\Sigma})} \quad (9)$$

Note that the explicit conditional dependence of $\boldsymbol{\theta}_r$ on \mathbf{z} has been suppressed, because it is understood that $\boldsymbol{\theta}_r$ corresponds to those components of $\boldsymbol{\theta}$ for which the corresponding z_i s take values of unity.

2.2.5. Sampling distribution of \mathbf{z}

In the derivation of sampling distribution of \mathbf{z} , the parameter $\boldsymbol{\theta}$ is marginalised to avoid dealing with the Dirac-delta functions; as a consequence, the conditional probability distributions of z_i s become dependent on \mathbf{y} . Additionally, the sampling distribution is marginalised over the noise-variance parameter σ^2 as this step leads to a quicker convergence of the Markov chains. As such, the resulting sampling distribution of \mathbf{z} , $p(\mathbf{z} | \mathbf{y}, p_0, v_s)$, is arrived at after performing integrations over both $\boldsymbol{\theta}$ and σ^2 :

$$p(\mathbf{z} | \mathbf{y}, p_0, v_s) \propto \left(\int \left(\int p(\mathbf{y}, \boldsymbol{\theta} | \mathbf{z}, v_s, \sigma^2) d\boldsymbol{\theta} \right) p(\sigma^2) d\sigma^2 \right) p(\mathbf{z} | p_0) \quad (10)$$

Since individual components of \mathbf{z} are independent of each other, the conditional probability distributions of z_i are derived separately. The computation of conditional probability distribution over z_i requires computing the probability of $z_i = 1$ compared to $z_i = 0$, given the same values of \mathbf{z}_{-i} , v_s and p_0 ; the term \mathbf{z}_{-i} denotes \mathbf{z} with its i^{th} component removed. Denote by ξ_i the probability with which one samples $z_i = 1$, then ξ_i can be expressed as:

$$\begin{aligned} \xi_i &= \frac{p(z_i = 1 | \mathbf{y}, \mathbf{z}_{-i}, v_s, p_0)}{p(z_i = 1 | \mathbf{y}, \mathbf{z}_{-i}, v_s, p_0) + p(z_i = 0 | \mathbf{y}, \mathbf{z}_{-i}, v_s, p_0)} \\ &= \frac{p(\mathbf{y} | z_i = 1, \mathbf{z}_{-i}, v_s) p(z_i = 1 | p_0)}{p(\mathbf{y} | z_i = 1, \mathbf{z}_{-i}, v_s) p(z_i = 1 | p_0) + p(\mathbf{y} | z_i = 0, \mathbf{z}_{-i}, v_s) p(z_i = 0 | p_0)} \\ &= \frac{p(\mathbf{y} | z_i = 1, \mathbf{z}_{-i}, v_s) p_0}{p(\mathbf{y} | z_i = 1, \mathbf{z}_{-i}, v_s) p_0 + p(\mathbf{y} | z_i = 0, \mathbf{z}_{-i}, v_s) (1 - p_0)} \\ &= \frac{p_0}{p_0 + \frac{p(\mathbf{y} | z_i = 0, \mathbf{z}_{-i}, v_s)}{p(\mathbf{y} | z_i = 1, \mathbf{z}_{-i}, v_s)} (1 - p_0)} = \frac{p_0}{p_0 + R_i (1 - p_0)} \quad \left[\text{where } R_i = \frac{p(\mathbf{y} | z_i = 0, \mathbf{z}_{-i}, v_s)}{p(\mathbf{y} | z_i = 1, \mathbf{z}_{-i}, v_s)} \right] \end{aligned} \quad (11)$$

As seen above, the evaluation of ξ_i requires the computation of the marginal likelihood $p(\mathbf{y} | \mathbf{z}, v_s)$, which will be derived next. The derivation of the marginal likelihood involves integrating the likelihood with respect to the distribution over $\boldsymbol{\theta}$, followed by that over σ^2 .

Marginalisation over $\boldsymbol{\theta}$

$$\begin{aligned} &p(\mathbf{y} | \mathbf{z}, v_s, \sigma^2) \\ &= \int p(\mathbf{y}, \boldsymbol{\theta} | \mathbf{z}, v_s, \sigma^2) d\boldsymbol{\theta} \\ &= \int p(\mathbf{y} | \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} | \mathbf{z}, v_s, \sigma^2) d\boldsymbol{\theta} \\ &= \int p(\mathbf{y} | \boldsymbol{\theta}_r, \sigma^2) p(\boldsymbol{\theta}_r | v_s, \sigma^2) d\boldsymbol{\theta}_r \quad [\because \text{Dirac-delta functions are integrated out}] \\ &= \int \mathcal{N}(\mathbf{y} | \mathbf{D}_r \boldsymbol{\theta}_r, \sigma^2 \mathbf{I}_N) \mathcal{N}(\boldsymbol{\theta}_r | \mathbf{0}, \sigma^2 v_s \mathbf{A}_{0,r}) d\boldsymbol{\theta}_r \\ &= \int \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{y} - \mathbf{D}_r \boldsymbol{\theta}_r)^T (\mathbf{y} - \mathbf{D}_r \boldsymbol{\theta}_r)}{2\sigma^2}\right) \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2}}{(2\pi v_s \sigma^2)^{r/2}} \exp\left(-\frac{\boldsymbol{\theta}_r^T \mathbf{A}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2 v_s}\right) d\boldsymbol{\theta}_r \\ &= \frac{1}{(2\pi\sigma^2)^{N/2}} \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2}}{(2\pi v_s \sigma^2)^{r/2}} \exp\left(-\frac{(\mathbf{y}^T \mathbf{y} - \mathbf{a}_N^T \mathbf{A}_N^{-1} \mathbf{a}_N)}{2\sigma^2}\right) \int \exp\left(-\frac{(\boldsymbol{\theta}_r - \mathbf{a}_N)^T \mathbf{A}_N^{-1} (\boldsymbol{\theta}_r - \mathbf{a}_N)}{2\sigma^2}\right) d\boldsymbol{\theta}_r \\ &= \frac{1}{(2\pi\sigma^2)^{N/2}} \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2} (\det(\mathbf{A}_N))^{1/2}}{(v_s)^{r/2}} \exp\left(-\frac{(\mathbf{y}^T \mathbf{y} - \mathbf{a}_N^T \mathbf{A}_N^{-1} \mathbf{a}_N)}{2\sigma^2}\right) \end{aligned}$$

where $\mathbf{A}_N = (\mathbf{D}_r^T \mathbf{D}_r + v_s^{-1} \mathbf{A}_{0,r}^{-1})^{-1}$ and $\mathbf{a}_N = \mathbf{A}_N \mathbf{D}_r^T \mathbf{y}$.

Marginalisation over σ^2

$$\begin{aligned}
& p(\mathbf{y} \mid \mathbf{z}, v_s) \\
&= \int p(\mathbf{y} \mid \mathbf{z}, v_s, \sigma^2) p(\sigma^2) d\sigma^2 \\
&= \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2} (\det(\mathbf{A}_N))^{1/2}}{(2\pi)^{N/2} (v_s)^{r/2}} \int \frac{1}{(\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{y}^T \mathbf{y} - \mathbf{a}_N^T \mathbf{A}_N^{-1} \mathbf{a}_N)}{2\sigma^2}\right) \mathcal{IG}(a_\sigma, b_\sigma) d\sigma^2 \\
&= \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2} (\det(\mathbf{A}_N))^{1/2}}{(2\pi)^{N/2} (v_s)^{r/2}} \frac{(b_\sigma)^{a_\sigma}}{\Gamma(a_\sigma)} \int \frac{1}{(\sigma^2)^{a_\sigma + N/2 + 1}} \exp\left(-\frac{b_\sigma + \frac{1}{2} (\mathbf{y}^T \mathbf{y} - \mathbf{a}_N^T \mathbf{A}_N^{-1} \mathbf{a}_N)}{\sigma^2}\right) d\sigma^2 \\
&= \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2} (\det(\mathbf{A}_N))^{1/2}}{(2\pi)^{N/2} (v_s)^{r/2}} \frac{(b_\sigma)^{a_\sigma}}{\Gamma(a_\sigma)} \frac{\Gamma(a_\sigma + \frac{N}{2})}{(b_\sigma + \frac{1}{2} (\mathbf{y}^T \mathbf{y} - \mathbf{a}_N^T \mathbf{A}_N^{-1} \mathbf{a}_N))^{(a_\sigma + \frac{N}{2})}}
\end{aligned}$$

where $\Gamma(\cdot)$ denotes the Gamma function. As such,

$$p(\mathbf{y} \mid \mathbf{z}, v_s) = \begin{cases} \frac{\Gamma(a_\sigma + \frac{N}{2})}{(2\pi)^{N/2}} \frac{(b_\sigma)^{a_\sigma}}{\Gamma(a_\sigma)} \frac{1}{(b_\sigma + \frac{1}{2} (\mathbf{y}^T \mathbf{y}))^{(a_\sigma + \frac{N}{2})}} & \text{when all } z_i \text{ are zero} \\ \frac{\Gamma(a_\sigma + \frac{N}{2})}{(2\pi)^{N/2} (v_s)^{r/2}} \frac{(b_\sigma)^{a_\sigma}}{\Gamma(a_\sigma)} \frac{(\det(\mathbf{A}_{0,r}^{-1}))^{1/2} (\det(\mathbf{A}_N))^{1/2}}{(b_\sigma + \frac{1}{2} (\mathbf{y}^T \mathbf{y} - \mathbf{a}_N^T \mathbf{A}_N^{-1} \mathbf{a}_N))^{(a_\sigma + \frac{N}{2})}} & \text{otherwise} \end{cases} \quad (12)$$

Finally, the sampling distributions of z_i , $i = 1, \dots, P$, can be computed using a Bernoulli distribution:

$$\boxed{z_i \mid \mathbf{y}, v_s, p_0 \sim \text{Bern}(\xi_i)} \quad (13)$$

with $\xi_i = \frac{p_0}{p_0 + R_i(1-p_0)}$ and $R_i = \frac{p(\mathbf{y} \mid z_i=0, \mathbf{z}_{-i}, v_s)}{p(\mathbf{y} \mid z_i=1, \mathbf{z}_{-i}, v_s)}$ calculated using Eq. (12). Lastly, it should be mentioned that to avoid numerical overflow errors, it is practical to use the logarithm of the marginal likelihoods $p(\mathbf{y} \mid \mathbf{z}, v_s)$ to compute $R_i = \exp\{\log[p(\mathbf{y} \mid z_i = 0, \mathbf{z}_{-i}, v_s)] - \log[p(\mathbf{y} \mid z_i = 1, \mathbf{z}_{-i}, v_s)]\}$.