

# Generalized additive models using Bayesian model selection with mixtures of g-priors

Gyeonghun Kang<sup>1</sup> and Seonghyun Jeong<sup>1,2</sup>

<sup>1</sup>Department of Statistics and Data Science, Yonsei University

<sup>2</sup>Department of Applied Statistics, Yonsei University

## Main Contributions

- Reviewed and extended estimation methods for GAM via BMS using g-priors.
- Proposed a simple slice sampler algorithm to draw from a class of generalized beta distributions, including truncated compound hypergeometric distribution (tCCH).
- Explained different model selection behaviors of mixtures of g-priors in terms of Bayes Factor of adding redundant variable.

## Bayesian Model Selection (BMS)

### Generalized Additive Model (GAM) via BMS

- GAM is an extension of a linear model assuming, for  $X_i \in R^p$ ,  $i = 1, \dots, n$ ,

$$h(E[y_i | X_i]) = \eta_i = \alpha + \sum_{j=1}^p f_j(x_{ij}), \quad \sum_{i=1}^n f_j(x_{ij}) = 0$$

- $f_j$  has a spline basis representation  $f_j(\cdot) = \sum_{k=1}^{K_j} \beta_{jk} b_{jk}(\cdot)$ , where the basis  $b_{jk}(\cdot)$  is linear for  $k = 1$  and the others are determined by the knots  $\xi_j = \{\xi_{j1}, \dots, \xi_{jL_j}\}$ .
- We use natural cubic spline basis for  $b_{jk}(\cdot)$  so that for  $K_j = 0$  it is reduced to linear.
- $\eta = (\eta_1, \dots, \eta_n)^T$  is written as  $\eta = \alpha 1_n + B\beta$  for  $B = [B_1, \dots, B_p] \in R^{n \times J}$ ,  $J = \sum_{j=1}^p K_j$ ,  $(B_j)_{i,k} = b_{jk}(x_{ij})$ , where  $B$  is column-wise centered for identifiability.
- We are interested in  $L: (\alpha, f_1, \dots, f_p) \mapsto L(\alpha, f_1, \dots, f_p)$  (e.g. pointwise estimate, CI, ...)

$$\pi(L(\alpha, f_1, \dots, f_p) | Y) = \int_{\Xi} \pi(L(\alpha, f_1, \dots, f_p) | \xi, Y) d\Pi(\xi | Y)$$

i.e., via Bayesian Model Selection. To this end, our Bayesian model formulation is

$$\begin{aligned} \Pi(\alpha, \beta_\xi, \xi) &= \Pi(\alpha) \Pi(\beta_\xi | \xi) \Pi(\xi) \\ \pi(\xi | Y) &\propto \pi(\xi) p(Y | \xi) \\ p(Y | \xi) &= \int \int p(Y | \alpha, \beta_\xi) d\Pi(\alpha) d\Pi(\beta_\xi | \xi) \end{aligned}$$

### Mixtures of g-priors for BMS

$$\frac{1}{g+1} \sim tCCH\left(\frac{a}{2}, \frac{b}{2}, \frac{s}{2}, \nu, \kappa\right)$$

$$\beta_\xi | g, \xi \sim N\left(0, g[\tilde{B}_\xi^T J_n(\hat{\eta}_\xi) \tilde{B}_\xi]^{-1}\right)$$

#### Mixtures of g prior

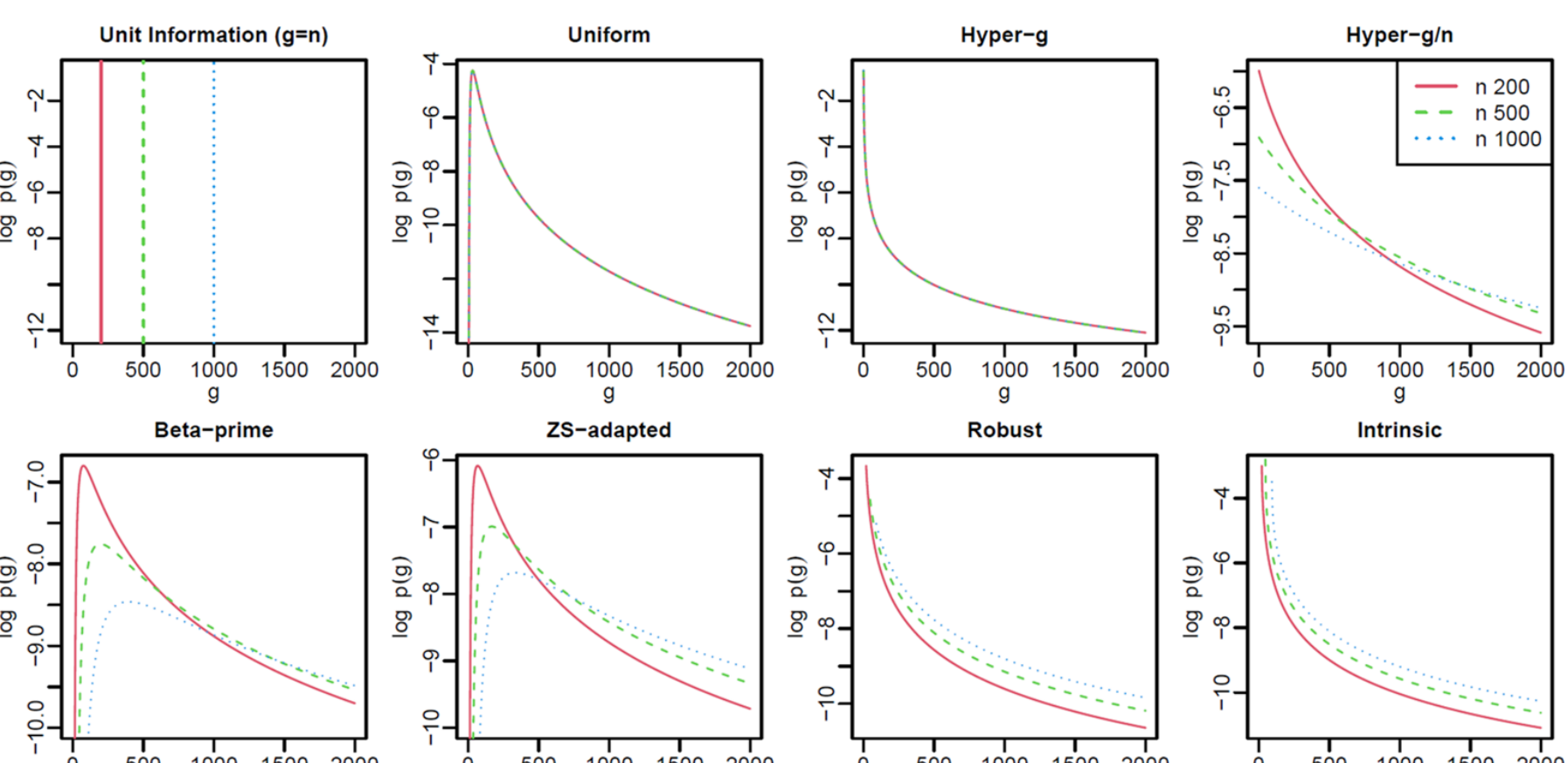
- For Laplace approximation of  $p(Y | \xi)$ , we use the variant of  $g$  prior for generalized linear model (Li and Clyde, 2018)

$$u \sim tCCH(a, b, z, s, \nu, \theta)$$

$$f(u) = \frac{\nu(\nu u)^{a-1}(1-\nu u)^{b-1}[\theta + (1-\theta)\nu u]^{-r} e^{-su}}{e^{-s/\nu} \Phi_1(b, r, a+b, s/\nu, 1-\theta) B(a, b)} \mathbf{1}_{\{0 < u < 1/\nu\}}$$

- Various mixtures of g-priors are classified into two groups according to the prior concentration:

$g = O(1)$	$g = O(n)$
Uniform, Hyper-g	Hyper-g/n Robust, Intrinsic, Beta-prime, ZS-adapted,

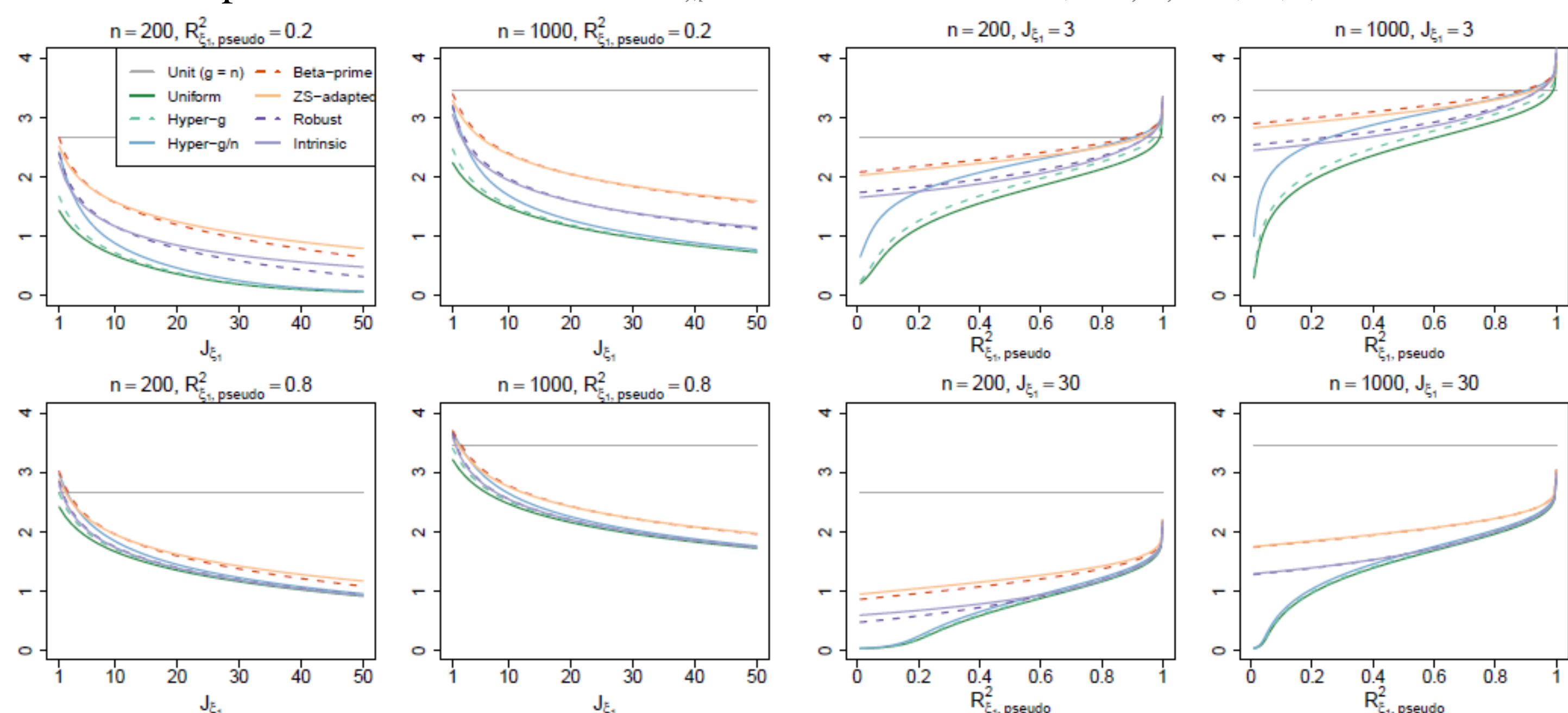


### Mixtures of g-priors as penalty functions

- The Bayes factor of two knots  $\xi_1, \xi_2$  where  $J_{\xi_1} = J_{\xi_2} + k$  ( $k \in N^+$ ) but  $\hat{\eta}_{\xi_1} = \hat{\eta}_{\xi_2}$  is

$$BF[\xi_1; \xi_2] = \begin{cases} (1+b)^{-k/2}, & \text{if } g \sim \delta_b(g), \\ E[(1+g)^{-k/2} | \xi_2, Y], & \text{if } g \sim tCCH(a/2, b/2, r, s/2, \nu, \kappa) \end{cases}$$

which is plotted below for  $k = 1$ :  $R_{\xi_1, \text{pseudo}}^2 = 1 - e^{-Q_{\xi_1}/n}$ , where  $Q_{\xi_1} = \beta_{\xi_1}^T \tilde{B}_{\xi_1}^T J_n(\hat{\eta}_{\xi_1}) \tilde{B}_{\xi_1} \beta_{\xi_1}$  is the Wald statistics.

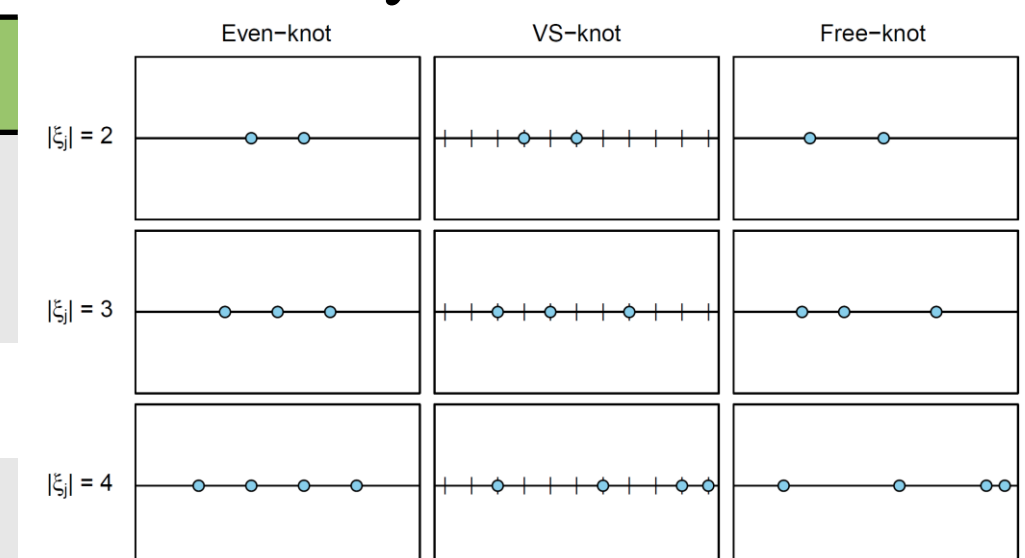


- $BF[\xi_1; \xi_2]$  is the penalty against the model  $\xi_1$  by allowing  $k$  redundant variables to no avail ( $\hat{\eta}_{\xi_1} = \hat{\eta}_{\xi_2}$ ). Whereas Unit information prior ( $g = n$ ) yields a constant penalty regardless of  $J_{\xi_1}$  and the goodness-of-fit,
  - mixtures of g-priors favor sparser models when comparing small models, but move towards more complex models when comparing large models, a trait desirable in capturing weak signals in the data.
  - The  $O(1)$  priors have the weakest penalty profiles among the mixtures of g-priors. Simulations showed that compared to the  $O(n)$ , the  $O(1)$  priors tend to overfit to noise in the data.

### Priors for knots

- A tradeoff exists between computational expedience and flexibility in estimates.

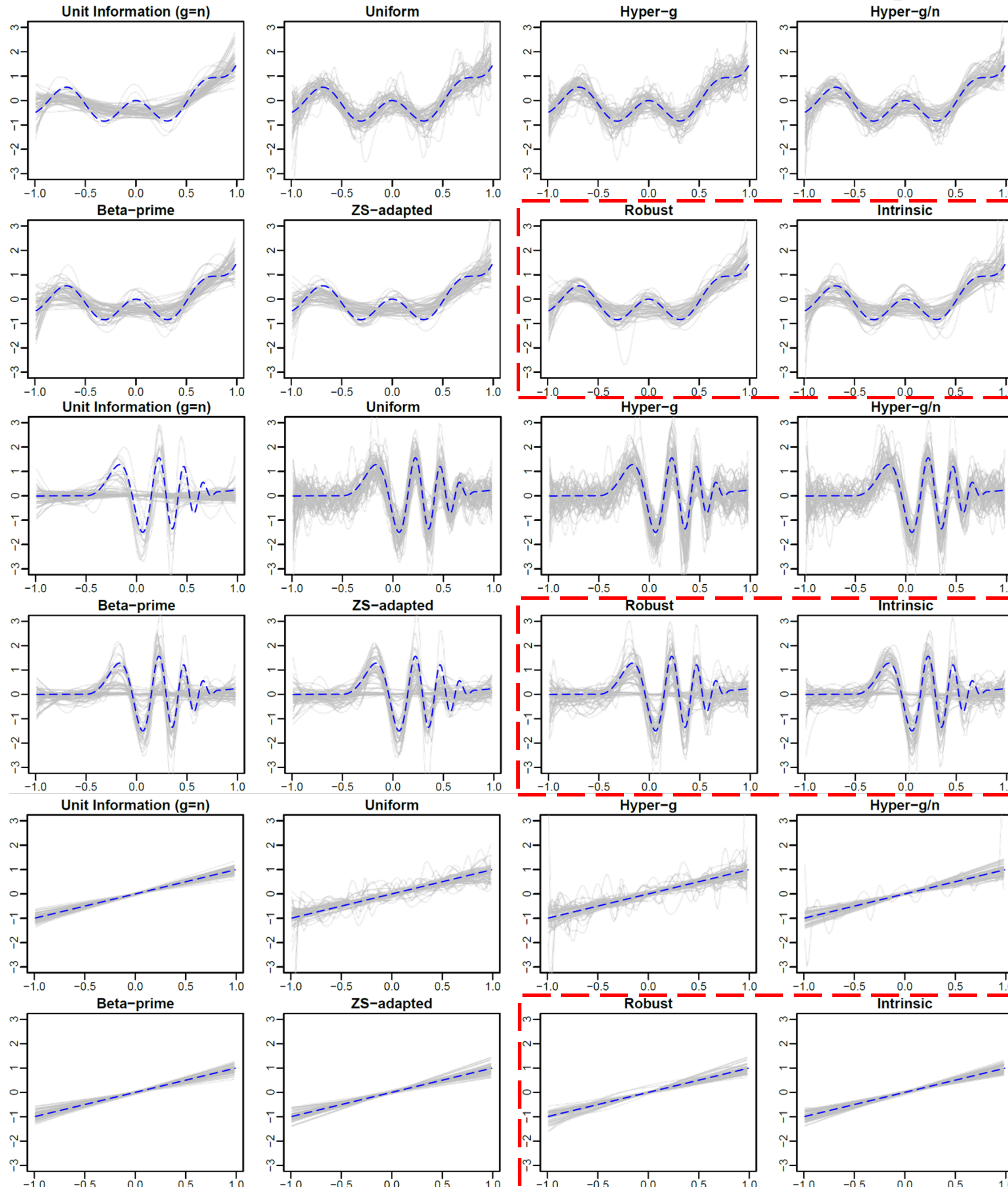
Knot Locations		Posterior Sampling
Even-knot	Equidistant points	Direct enumeration
		Metropolis Hastings
VS-knot	Among grid points	Gibbs Sampling
Free-knot	Anywhere	Metropolis Hastings



## Simulation Results

### Comparison among the mixtures of g-priors

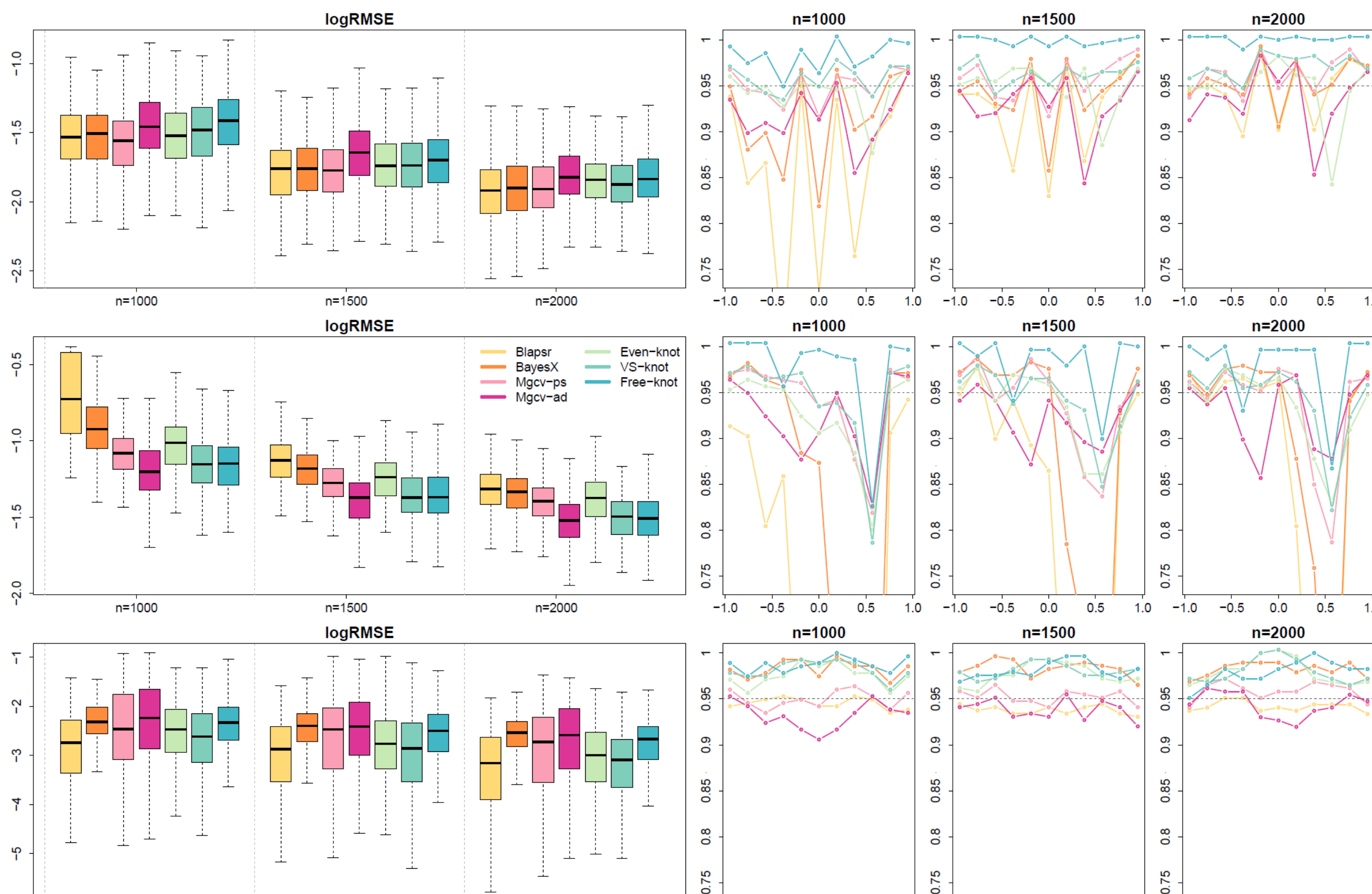
- For  $f_j, j = 1, 2, 3$ , test functions, we generated the univariate model  $\eta_i = \alpha + f_j(x_i)$  where  $x_j \sim \text{Unif}(-1, 1)$  and the response follows Bernoulli, Poisson and Gaussian, for  $n = 100, 200, 300$  ( $n = 500, 1000, 2000$  for Bernoulli) samples. We used VS-knot.



The result for  $n = 500$  Bernoulli univariate model fitted via VS-knot. The true function is plotted in blue dashed and the posterior means for 50 of 500 replications in grey lines.

### Comparison with other methods

- The competitors, all based on the idea of Bayesian P-splines, include **R2BayesX**, **Blapsr**, and **Mgcv** with locally adaptive (**Mgcv-ad**) and non-adaptive estimation (**Mgcv-ps**).
- For  $f_j, j = 1, 2, 3$ , test functions, we generated the univariate model  $\eta_i = \alpha + \sum_{j=1}^3 f_j(x_i)$  where  $x_j \sim \text{Unif}(-1, 1)$  and the response follows Bernoulli, Poisson and Gaussian, for  $n = 100, 200, 300$  ( $n = 1000, 1500, 2000$  for Bernoulli) samples.
- R2BayesX** and **Blapsr** often oversmooth the targets with excessive penalization, while **Mgcv** provides too wiggly estimates of the linear function, implying undersmoothing. **VS-knot** outperforms in most cases in terms of logRMSE and coverage probability.



The log RMSE and the 95% coverage prob. in the logistic models with  $n = 1000, 1500, 2000$ , for  $j = 1, 2, 3$  from top to bottom.

## Real Data Applications (Pima Diabetes Data)

- Signs of diabetes (binary) of  $n = 532$  women in Pima Indian population, Arizona.
- Our BMS-based methods provide posterior probability that a function is indeed linear.

$$\text{logit}(p_i) = \alpha + f_1(\text{pregnant}_i) + f_2(\text{glucose}_i) + f_3(\text{pressure}_i) + f_4(\text{mass}_i) + f_5(\text{pedigree}_i) + f_6(\text{age}_i)$$

