

Power-Expected-Posterior Priors in Generalized Linear Models

Ioannis Ntzoufras,

*Department of Statistics, Athens University of Economics and Business, Athens, Greece; e-mail:
ntzoufras@aueb.gr.*

Joint work with:

Dimitris Fouskakis & Konstantinos Perrakis

Department of Mathematics

National Technical University of Athens

Department of Statistics

Athens University of Economics and Business

Available at <http://stat-athens.aueb.gr/~jbn/papers/obayes15.htm>.

Synopsis

1. From the expected-posterior prior (EPP) to the power-expected-posterior (PEP) prior
2. Alternative definitions of the power likelihood in PEP-priors
3. Implementing the method in GLMs (MCMC algorithm)
4. Illustrations
5. Using Mixtures of PEP priors
6. Illustrations (continued)
7. Discussion

1 Introduction: Model Selection and Expected-Posterior Priors

Within the Bayesian framework the comparison between models M_0 and M_1 is evaluated via the **Posterior Odds** (PO)

$$PO_{01} \equiv \frac{\pi(M_0|\mathbf{y})}{\pi(M_1|\mathbf{y})} = \frac{m_0(\mathbf{y})}{m_1(\mathbf{y})} \times \frac{\pi(M_0)}{\pi(M_1)} = BF_{01} \times O_{01} \quad (1)$$

which is a function of the **Bayes Factor** (BF_{01}) and the **Prior Odds** (O_{01}).

In the above $m_\ell(\mathbf{y})$ is the marginal likelihood under model M_ℓ and $\pi(M_\ell)$ is the prior probability of model M_ℓ .

The marginal likelihood of model M_ℓ is given by

$$m_\ell(\mathbf{y}) = \int f_\ell(\mathbf{y}|\boldsymbol{\theta}_\ell)\pi_\ell(\boldsymbol{\theta}_\ell)d\boldsymbol{\theta}_\ell, \quad (2)$$

where $f_\ell(\mathbf{y}|\boldsymbol{\theta}_\ell)$ is the likelihood under model M_ℓ with parameters $\boldsymbol{\theta}_\ell$ and $\pi_\ell(\boldsymbol{\theta}_\ell)$ is the prior distribution of model parameters given model M_ℓ .

Expected-Posterior Priors (EPP)

- Pérez & Berger (2002, Biometrika) developed priors for use in model comparison, through utilization of the device of **imaginary training samples**.
- They defined the **expected-posterior prior** (EPP) as the posterior distribution of a parameter vector for the model under consideration, averaged over all possible imaginary samples $\mathbf{y}^* = (y_1^*, \dots, y_{n^*}^*)^T$ coming from a “suitable” predictive distribution $m^*(\mathbf{y}^*)$.

Hence the EPP for the parameters of any model M_ℓ is

$$\pi_\ell^{EPP}(\boldsymbol{\theta}_\ell) = \int \pi_\ell^N(\boldsymbol{\theta}_\ell | \mathbf{y}^*) m^*(\mathbf{y}^*) d\mathbf{y}^*, \quad (3)$$

where $\pi_\ell^N(\boldsymbol{\theta}_\ell | \mathbf{y}^*)$ is the posterior of $\boldsymbol{\theta}_\ell$ for model M_ℓ using a baseline prior $\pi_\ell^N(\boldsymbol{\theta}_\ell)$ and data \mathbf{y}^* .

Features of EPP

- **Impropriety** of baseline priors causes no indeterminacy. **Impropriety** in m^* also does not cause indeterminacy, because it is common to the EPPs for all models.
- It makes priors **compatible** across models, through their dependence on a common data distribution.
- Usually we consider as m^* the marginal likelihood of a reference model.
 - Usual choices in regression models are the **null** and full **model**.
 - Here we consider the null, i.e. $m^*(\mathbf{y}^*) = m_0^N(\mathbf{y}^*)$.
- In nested cases usually the reference model is the simplest model. In this case EPP is the same as the **Intrinsic Prior**.
- We choose the smallest n^* for which the posterior is proper: **minimal training sample size**.
- **Main Issue**: In variable selection problems specification of X_ℓ^* .

Power-Expected-Posterior (PEP) Priors

Fouskakis, Ntzoufras and Draper (2015, *Bayesian Analysis*).

$$\begin{array}{c}
 \underbrace{\pi_{\ell}^{EPP}(\boldsymbol{\theta}_{\ell})}_{\Downarrow} = \int \underbrace{\pi_{\ell}^N(\boldsymbol{\theta}_{\ell}|\mathbf{y}^*)}_{\Downarrow} \underbrace{m_0^N(\mathbf{y}^*)}_{\Downarrow} d\mathbf{y}^* \\
 \pi_{\ell}^{PEP}(\boldsymbol{\theta}_{\ell}; \delta) = \int \underbrace{\pi_{\ell}^N(\boldsymbol{\theta}_{\ell}|\mathbf{y}^*, \delta)}_{\Downarrow} \underbrace{m_0^N(\mathbf{y}^*|\delta)}_{\Downarrow} d\mathbf{y}^*
 \end{array}$$

we substitute the likelihood terms with powered-versions of the likelihoods
(i.e. they are raised to the power of $1/\delta$).

Features of PEP

PEP priors method amalgamates ideas from Intrinsic Priors, EPPs, Unit Information Priors and Power Priors, to unify ideas of Non-Data Objective Priors.

PEP priors solve the following problems:

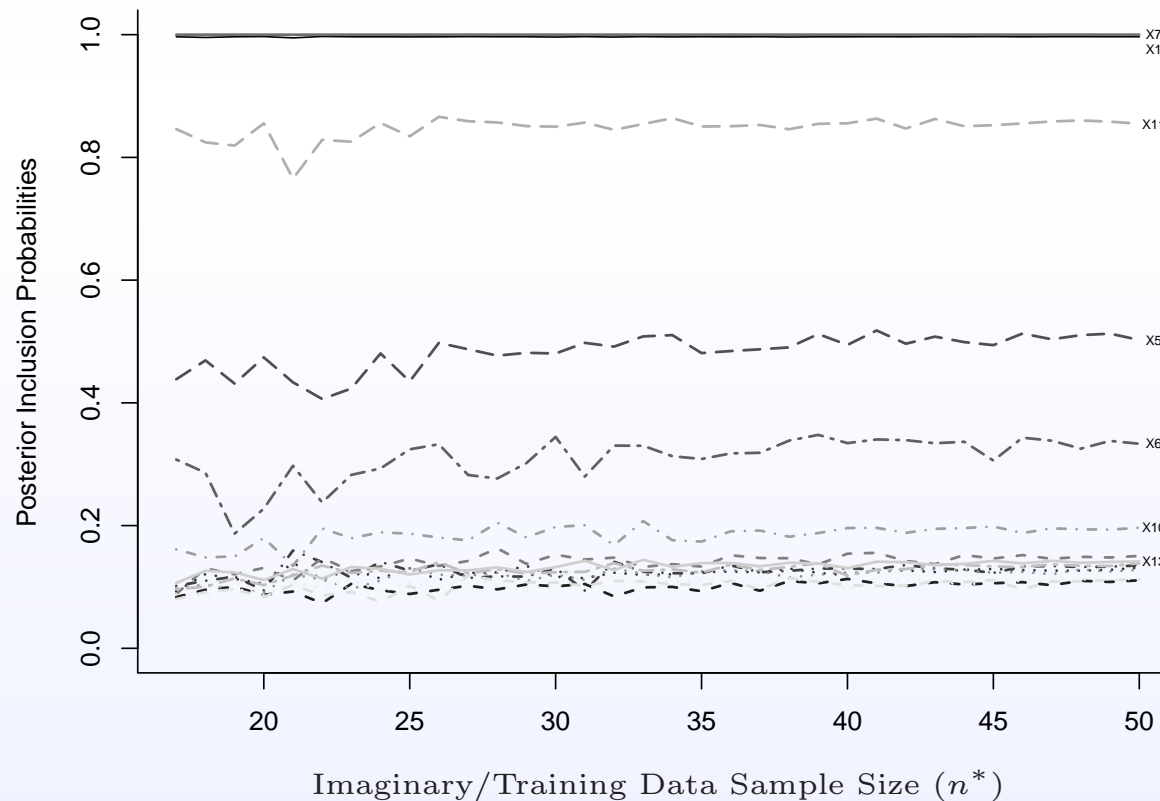
- Dependence of training sample size.
- Lack of robustness with respect to the sample irregularities.
- Excessive weight of the prior when the number of parameters is close to the number of data.

At the same time the PEP prior is a fully objective method and shares the advantages of Intrinsic Priors and EPPs.

- We choose $\delta = n^*$, $n^* = n$ and therefore $X_\ell^* = X_\ell$; by this way we dispense with the selection of the training samples.

Sensitivity analysis on imaginary sample size

Figure 1: *Posterior marginal inclusion probabilities, for n^* values from 17 to $n = 50$, with the PEP prior methodology (simulated example for a variable selection problem in normal linear model).*



Features of PEP (cont.)

For Normal models

- In Fouskakis, Ntzoufras & Draper, 2015 (*Bayesian Analysis*) we illustrated the the PEP prior approach
 - is robust with respect to the training sample size
 - is not informative when d_ℓ is close to n .
- The PEP prior can be expressed as a mixture of g -priors (Fouskakis, Ntzoufras & Pericchi, *unpublished work, presented in ISBA2014*).
- The Power-conditional-expected-posterior (PCEP) prior (Fouskakis & Ntzoufras, 2015, *to appear in JCGS*) is similar to the g -prior with (i) more complicated variance structure, (ii) more dispersed and (iii) more parsimonious than the g -prior
- Both PEP and PCEP are leading to consistent variable selection methods.

2 Extension to Generalized Linear Models

Definitions of the power-likelihood

Normal regression models: the definition of the power-likelihood seems quite clear.

We have worked with the density-normalized power likelihood since for any normal distribution with mean μ and variance σ^2 it holds that

$$f(y|\mu, \sigma^2, \delta) = \frac{f(y|\mu, \sigma^2)^{1/\delta}}{\int f(y|\mu, \sigma^2)^{1/\delta} dy} = N(\mu, \delta \sigma^2)$$

This is not the case for all distributions in the exponential family and hence for GLMs.

Definitions of the power-likelihood

Density-normalized power likelihoods in GLMs: May end up to a distribution which is not the same as the one in the original model formulation.

- In binary logistic regression \Rightarrow power likelihood is still Bernoulli with success probability

$$\frac{\pi^{1/\delta}}{\pi^{1/\delta} + (1 - \pi)^{1/\delta}}.$$

- This is not the case for the Binomial and the Poisson models resulting is some cumbersome distributions which increase computational complexity (without any obvious gain).

Alternative definitions of the power-likelihood

We consider the PEP representation

$$\pi_{\ell}^{PEP}(\boldsymbol{\theta}_{\ell}; \delta) = \int \pi_{\ell}^N(\boldsymbol{\theta}_{\ell} | \mathbf{y}^*, \delta) m_0^N(\mathbf{y}^* | \delta) d\mathbf{y}^*$$

with δ controlling the amount of prior-information accounted in the final posterior (and the dispersion of the prior distribution).

We now consider **the unnormalized power-likelihood** and then normalize the posterior (which is also the approach in Ibrahim and Chen, 2000, *Stat.Science*).

Hence

$$\pi_{\ell}^N(\boldsymbol{\theta}_{\ell} | \mathbf{y}^*, \delta) = \frac{f_{\ell}(\mathbf{y}^* | \boldsymbol{\theta}_{\ell})^{1/\delta} \pi_{\ell}^N(\boldsymbol{\theta}_{\ell})}{\int f_{\ell}(\mathbf{y}^* | \boldsymbol{\theta}_{\ell})^{1/\delta} \pi_{\ell}^N(\boldsymbol{\theta}_{\ell}) d\boldsymbol{\theta}_{\ell}}$$

What about $m_0^N(\mathbf{y}^* | \delta)$?

Two alternatives for the marginal distribution

- Consider the **unnormalized power-likelihood** and then normalize m_0^N :

$$m_0^N(\mathbf{y}^*, \delta) = \frac{\int f_0(\mathbf{y}^* | \boldsymbol{\theta}_0)^{1/\delta} \pi_0^N(\boldsymbol{\theta}_0) d\boldsymbol{\theta}_0}{\int \int f_0(\mathbf{y}^* | \boldsymbol{\theta}_0)^{1/\delta} \pi_0^N(\boldsymbol{\theta}_0) d\boldsymbol{\theta}_0 d\mathbf{y}^*} .$$

This will be noted as the *Diffuse Reference PEP (DR-PEP)*.

- Consider the **original likelihood** (without introducing any further uncertainty) i.e.

$$m_0^N(\mathbf{y}^*, \delta) = m_0(\mathbf{y}^*) = \int f_0(\mathbf{y}^* | \boldsymbol{\theta}_0) \pi_0^N(\boldsymbol{\theta}_0) d\boldsymbol{\theta}_0 .$$

This will be noted as the *Concentrated Reference PEP (CR-PEP)*.

In both cases the expected-posterior interpretation is retained with the first prior being more diffuse than the second.

Features of the diffuse-reference PEP

- Still has the interpretation of a posterior density given some imaginary data \mathbf{y}^* “weighted” by n^*/δ data-points and averaged over a data distribution.
- The same type of uncertainty is introduced both in the “posterior” and the predictive (averaged) part.
- In normal regression models
 - Equivalent to using the the density-normalized power likelihood.
 - It is equivalent to PEP and PCEP.
 - It leads to a consistent model selection method.
 - It is more dispersed (and parsimonious) than the g-prior.

Features of the concentrated-reference PEP

- Still has the interpretation of a posterior density given some imaginary data \mathbf{y}^* “weighted” by n^*/δ data-points averaged over the predictive distribution of **the actual reference model**.
- Different type of uncertainty is introduced both in the “posterior” and the predictive (averaged) part.
- Less dispersed than the diffuse version of PEP.
- In normal regression models
 - It is less dispersed (and parsimonious) than PEP (and DR-PEP) and more dispersed (and parsimonious) than the g-prior.
 - It leads to a consistent model selection method.

Comparison of the two approaches in normal regression

Volume variance multipliers in normal regression models

The volume of the variance-covariance matrix in the g-prior and in the two PEP approaches is given by

$$\left| \text{Var}(\boldsymbol{\beta}_\ell | M_\ell) \right| = \varphi(n, d_\ell) \times |\mathbf{X}_\ell^T \mathbf{X}_\ell|^{-1}$$

- **G-prior** with $g = n \Rightarrow \varphi(n, d_\ell) = n^{d_\ell}$
- **DR-PEP prior** $\Rightarrow \varphi(n, d_\ell) = n^{2d_\ell} \left[\frac{2n+1}{(n+1)^2} \right]^{d_\ell - d_0}$
- **CR-PEP prior** $\Rightarrow \varphi(n, d_\ell) = n^{d_\ell} \left[\frac{n^2+2n}{n^2+2n+1} \right]^{d_\ell} \left[\frac{n^2+n+2}{n+2} \right]^{d_0}$.

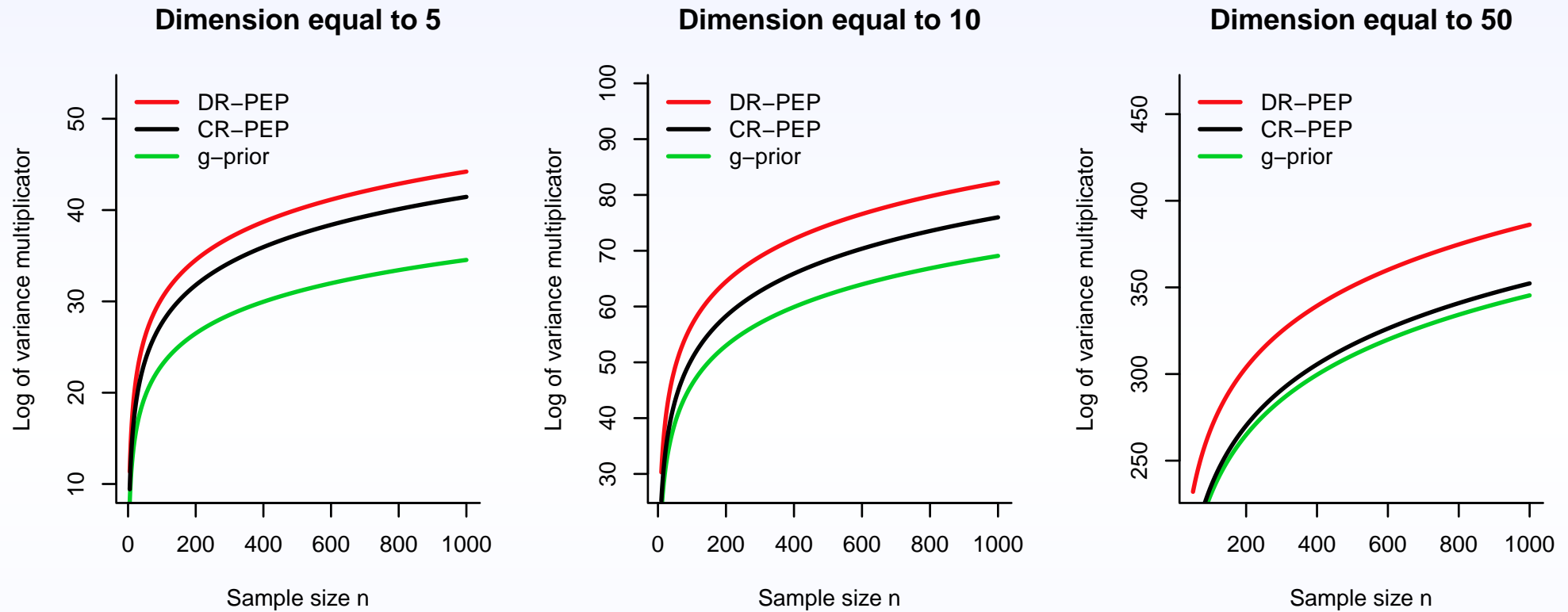


Figure 2: Log-variance multipliers of the DR-PEP, CR-PEP and g -priors versus sample size for $d_\ell = 5, 10, 50$.

Formulation in GLMs

- $M_\ell \rightarrow \gamma$: Binary variable inclusion indicators (γ) in order to search the model space using Gibbs sampling (George and McCulloch, 1993, *JASA*)
- $Y_i \sim$ a distribution member of the exponential family.
The parameters of the distribution are associated with the linear predictor via a link function.
- p covariates.
- \mathbf{X} is the $n \times (p + 1)$ data matrix with the first column to be the constant and rest containing the data of each covariate.
- \mathbf{X}_γ is the $n \times d_\gamma$ data matrix for model γ with $d_\gamma = \sum_{j=0}^p \gamma_j$ covariates.
- β_γ is the parameter vector of length d_γ with the effects of each covariate
- The linear predictor vector is given by $\eta_\gamma = \mathbf{X}_\gamma \beta_\gamma$

We focus our presentation on the DR-PEP (computation is similar for the CR-PEP)

The prior

$$\pi_{\gamma}^{DRPEP}(\boldsymbol{\beta}_{\gamma}) \propto \int \int \left\{ \frac{f_{\gamma}(\mathbf{y}^* | \boldsymbol{\beta}_{\gamma})^{1/\delta} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma})}{\int f_{\gamma}(\mathbf{y}^* | \boldsymbol{\beta}_{\gamma})^{1/\delta} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma}) d\boldsymbol{\beta}_{\gamma}} \right\} f_0(\mathbf{y}^* | \boldsymbol{\beta}_0)^{1/\delta} \pi_0^N(\boldsymbol{\beta}_0) d\boldsymbol{\beta}_0 d\mathbf{y}^*$$

Two possible approaches to simplify the above expression

- The **posterior part** can be well approximated by a normal distribution (Chen and Ibrahim, 2003, *Stat.Sinica*)
- Integral in the denominator can be well approximated using Laplace approximation

The posterior

$$\pi_{\gamma}^{DRPEP}(\boldsymbol{\beta}_{\gamma}|\mathbf{y}) \propto f_{\gamma}(\mathbf{y}|\boldsymbol{\beta}_{\gamma}) \int \int \frac{f_{\gamma}(\mathbf{y}^*|\boldsymbol{\beta}_{\gamma})^{1/\delta} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma})}{\int f_{\gamma}(\mathbf{y}^*|\boldsymbol{\beta}_{\gamma})^{1/\delta} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma}) d\boldsymbol{\beta}_{\gamma}} f_0(\mathbf{y}^*|\boldsymbol{\beta}_0)^{1/\delta} \pi_0^N(\boldsymbol{\beta}_0) d\boldsymbol{\beta}_0 d\mathbf{y}^*$$

The marginal likelihood

$$m_{\gamma}^{DRPEP}(\mathbf{y}) \propto \int \int \int f_{\gamma}(\mathbf{y}|\boldsymbol{\beta}_{\gamma}) \frac{f_{\gamma}(\mathbf{y}^*|\boldsymbol{\beta}_{\gamma})^{1/\delta} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma})}{\int f_{\gamma}(\mathbf{y}^*|\boldsymbol{\beta}_{\gamma})^{1/\delta} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma}) d\boldsymbol{\beta}_{\gamma}} f_0(\mathbf{y}^*|\boldsymbol{\beta}_0)^{1/\delta} \pi_0^N(\boldsymbol{\beta}_0) d\boldsymbol{\beta}_{\gamma} d\boldsymbol{\beta}_0 d\mathbf{y}^*$$

In order to estimate the posterior model probabilities, we use an MCMC scheme with full data augmentation by introducing

- For each model γ , we introduce a complement of β_γ denoted by $\beta_{\setminus\gamma}$ for all coefficients not included in the model.
- A pseudoprior $\pi_\gamma(\beta_{\setminus\gamma})$ is defined to play the role of a proposal and the linear predictor can be rewritten as $\eta_i = \sum_{j=0}^p X_{ij} \gamma_j b_{\gamma,j}$ where $b_{\gamma,j}$ is the element of $\mathbf{b}_\gamma = (\beta_\gamma, \beta_{\setminus\gamma})$ which corresponds to covariate X_j .
- A latent parameter β_0 for the parameter of the reference model
- A latent vector of imaginary data \mathbf{y}^*

- We build a Gibbs based variable selection algorithm providing samples from the augmented posterior

$$\pi_{\gamma}^{DRPEP}(\boldsymbol{\beta}_{\gamma}, \boldsymbol{\beta}_{\setminus\gamma}, \gamma, \mathbf{y}^*, \beta_0 | \mathbf{y})$$

$$\propto \frac{f_{\gamma}(\mathbf{y} | \boldsymbol{\beta}_{\gamma}) \left[f_{\gamma}(\mathbf{y}^* | \boldsymbol{\beta}_{\gamma}) f_0(\mathbf{y}^* | \beta_0) \right]^{1/\delta}}{\int f_{\gamma}(\mathbf{y}^* | \boldsymbol{\beta}_{\gamma})^{1/\delta} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma}) d\boldsymbol{\beta}_{\gamma}} \pi_{\gamma}^N(\boldsymbol{\beta}_{\gamma}) \pi_{\gamma}^N(\boldsymbol{\beta}_{\setminus\gamma}) \pi_0^N(\beta_0) \pi(\gamma)$$

- We use Laplace approximation to evaluate the integral in the denominator.
- In this work, we use the Jeffreys prior as a baseline prior.

The MCMC algorithm - Gibbs variable selection for PEP

For each iteration t ($t = 1, 2, \dots, N$),

Step 1: For $j = 1, \dots, p$, we update $\gamma_j \sim \text{Bernoulli}\left(\frac{O_j}{1+O_j}\right)$, with

Step 2: We update β_γ from the full conditional posterior (given the current values of γ and \mathbf{y}^*) using a Metropolis step and proposals build using MLEs from a model with response $(\mathbf{y}, \mathbf{y}^*)$ and weights $\mathbf{w} = (\mathbf{1}_n, \delta^{-1}\mathbf{1}_{n^*})$

Step 3: Update $\beta_{\setminus\gamma}$ from the pseudo-prior $\pi_\gamma(\beta_{\setminus\gamma}) = N_{d_{\setminus\gamma}}\left(\hat{\beta}_{\setminus\gamma}, \mathbf{I}_{d_{\setminus\gamma}} \hat{\sigma}_{\beta_{\setminus\gamma}}^2\right)$.

Step 4: Sample β_0 from the full conditional posterior (given \mathbf{y}^*) using a Metropolis step with a normal proposal with mean the MLE with response \mathbf{y}^* and variance equal to $\delta \hat{\sigma}_{\hat{\beta}_0}^2$ with the latter being the corresponding variance of the MLE.

Step 5: Sample \mathbf{y}^* from the full conditional posterior (given β_γ , β_0 and γ) using a Metropolis step.

- The proposal depends on the model likelihood i.e. the stochastic part of the model; for details see next slide.
- In the acceptance probabilities we need to evaluate the marginal likelihoods $m_\gamma^N(\mathbf{y}^*|\delta)$ and $m_\gamma^N(\mathbf{y}^{*'}|\delta)$ which are computed by using Laplace approximation.

Details about the proposal for y_i^*

For \mathbf{y}^* we construct proposals depending on the likelihood of the model.

Binomial response: A product binomial proposal distribution is used with probability of success equal to

$$\pi_i = \frac{(\pi_{0,i}^D \pi_{\gamma,i})^{1/\delta}}{(\pi_{0,i}^D \pi_{\gamma,i})^{1/\delta} + [(1 - \pi_{0,i})^D (1 - \pi_{\gamma,i})]^{1/\delta}}$$

with $D = 1$ for DR-PEP and $D = \delta$ for CR-PEP, $\pi_{0,i} = [1 + \exp(-\beta_0)]^{-1}$ and $\pi_{\gamma,i} = [1 + \exp(-\eta_{\gamma,i})]^{-1}$; where $\eta_{\gamma,i}$ is the i -th element of $\boldsymbol{\eta}_{\gamma} = \mathbf{X}_{\gamma} \boldsymbol{\beta}_{\gamma}$.

Poisson regression:

CR-PEP: A product Poisson proposal distribution is used with mean equal to $\lambda_i = \lambda_{0,i} \lambda_{\gamma,i}^{1/\delta}$, with $\lambda_{0,i} = \exp(\beta_0)$ and $\lambda_{\gamma,i} = \exp(\eta_{\gamma,i})$.

DR-PEP: The same strategy for DR-PEP failed and we have used a simple Poisson proposal with mean $\lambda_i = y_i^*$.

Further Remarks about the MCMC

- Metropolis-Hastings steps are not needed for $\beta_{\setminus\gamma}$ and γ
 - $\beta_{\setminus\gamma}$ is sampled directly from the pseudo-prior distribution.
 - γ is sampled directly from the full conditional Bernoulli distribution.
- The pseudo-prior of $\beta_{\setminus\gamma}$ serves the role of the proposal and it does not influence the posterior but it does influence the efficiency of the MCMC algorithm.
- No specific fine tuning is required for the proposal distributions of $\beta_{\setminus\gamma}$ and β_0 (normal proposals based on MLEs).

Illustrative example 1: Simulated Binomial data

$n = 200$ binary responses with $p = 10$ potential covariates.

For $i = 1, \dots, n$

$$X_{ij} \sim N(0, 1) \text{ for } j = 1, \dots, 5$$

$$X_{ij} \sim N(0.3X_{i1} + 0.5X_{i2} + 0.7X_{i3} + 0.9X_{i4} + 1.1X_{i5}, 1) \text{ for } j = 6, \dots, 10$$

$$Y_i \sim \text{Bernoulli}(p_i)$$

Three scenarios:

Null: $\text{logit}(p_i) = 0.1$

Sparse: $\text{logit}(p_i) = 0.1 - 0.9X_{i3} + 1.2X_{i7} + 0.4X_{i10}$

Dense: $\text{logit}(p_i) = 0.1 + 0.6X_{i1} - 0.9X_{i3} + X_{i5} + 0.9X_{i6} + 1.2X_{i7} - 1.2X_{i8} - 0.5X_{i9}$

Scenario 1 (Null): True model = null

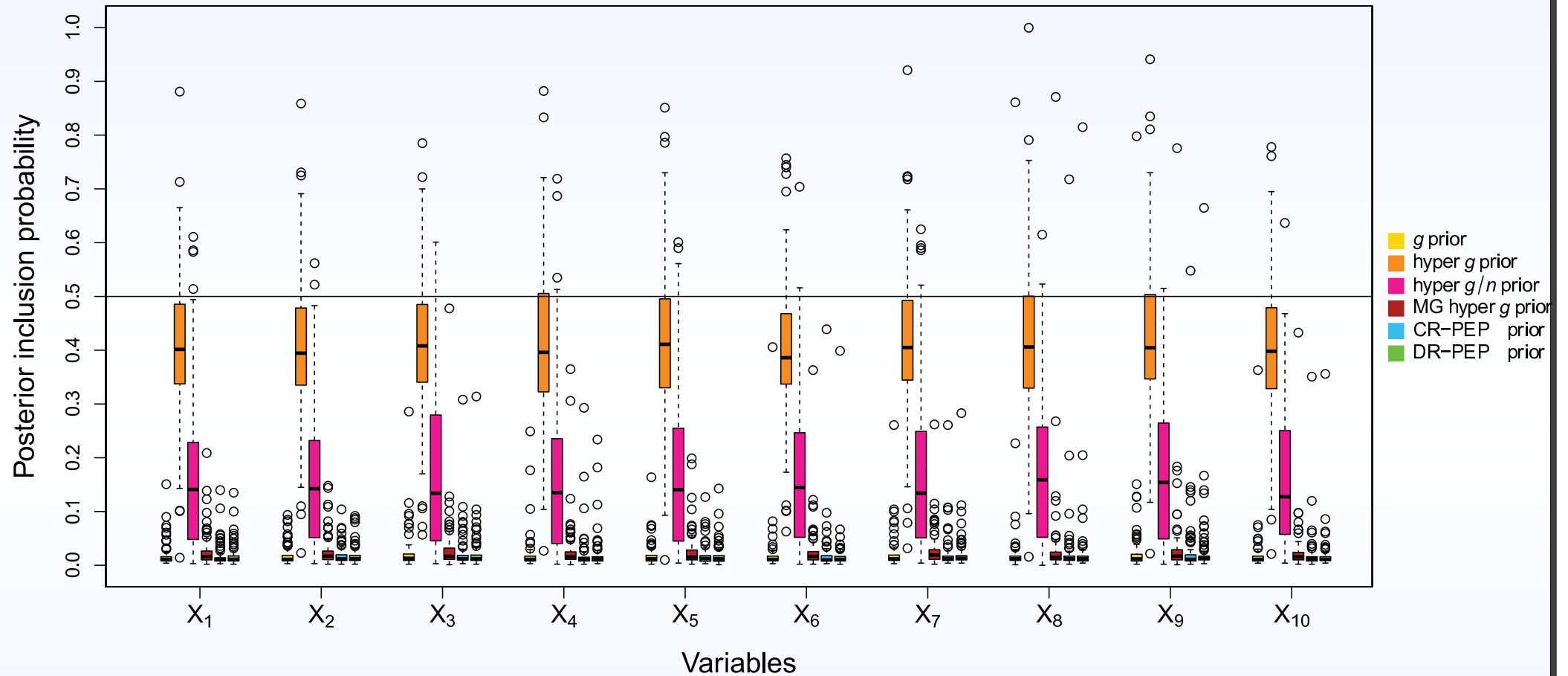


Figure 3: Posterior inclusion probabilities from 100 samples.

MG hyper- g prior: Maruyama & George (2011, Annals of Statistics) prior.

Scenario 2 (Sparse): True model = $0.1 - 0.9X_3 + 1.2X_7 + 0.4X_{10}$

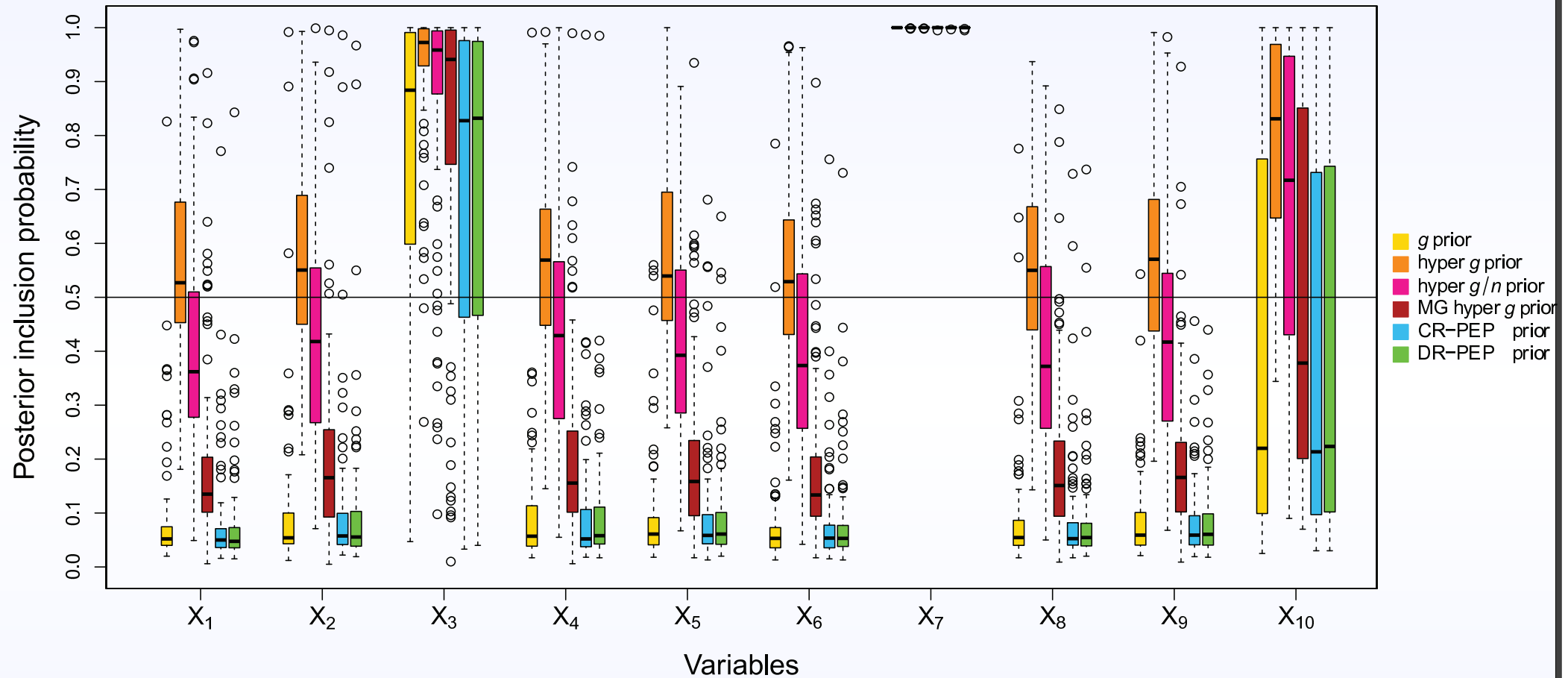


Figure 4: Posterior inclusion probabilities from 100 samples.

MG hyper-g prior: Maruyama & George (2011, Annals of Statistics) prior.

Scenario 3 (Dense)

True model = $0.1 + 0.6X_1 - 0.9X_3 + X_5 + 0.9X_6 + 1.2X_7 - 1.2X_8 - 0.5X_9$

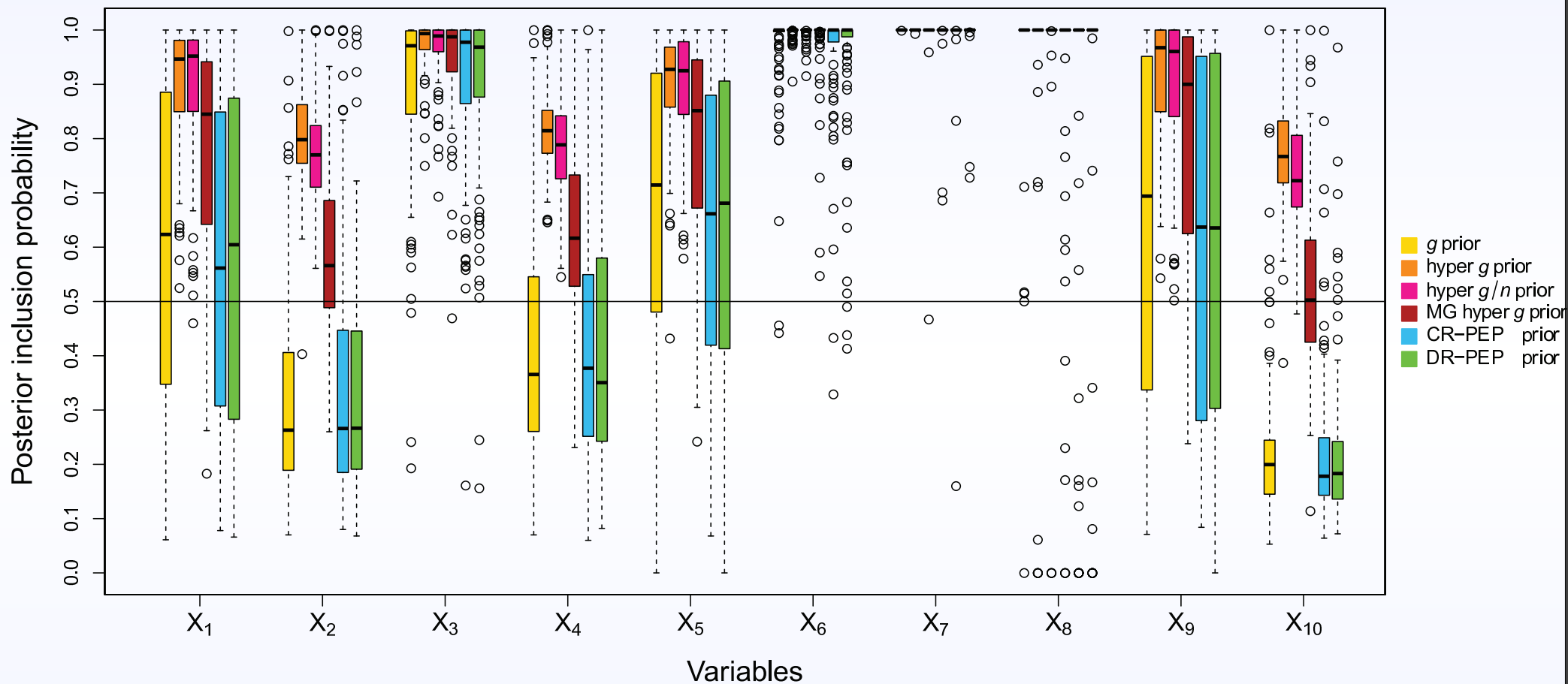


Figure 5: Posterior inclusion probabilities from 100 samples.

MG hyper-g prior: Maruyama & George (2011, Annals of Statistics) prior.

3 Hyper-delta PEP priors

PEP priors with fixed δ are similar in notion and behaviour as the g-priors.

We extend our approach by using hyper-priors for δ in a similar manner as hyper-g priors do.

Under this setting, the hyper- δ PEP prior can be approximated by

$$\pi_{\gamma}^{\text{PEP}}(\boldsymbol{\beta}_{\gamma}) \approx \int \int f_{N_{d_{\gamma}}}(\boldsymbol{\beta}_{\gamma}; \widehat{\boldsymbol{\beta}}_{\gamma}^*, \delta (\mathbf{X}_{\gamma}^{*T} \mathbf{H}_{\gamma}^* \mathbf{X}_{\gamma}^*)^{-1}) m_0^N(\mathbf{y}^* | \delta) \pi(\delta) d\mathbf{y}^* d\delta, \quad (4)$$

where $\widehat{\boldsymbol{\beta}}_{\gamma}^*$ is the MLE given the imaginary data.

This approximation cannot be applied when using EPPs with minimal training samples.

Similarly to the hyper-g (Liang *et al.*, 2008, *JASA*), the hyper-delta prior is given by

$$\pi(\delta) = \frac{a-2}{2} (1+\delta)^{-a/2},$$

which introduces the following prior for $\delta/(1+\delta)$

$$\frac{\delta}{1+\delta} \sim \text{Beta}\left(1, \frac{a}{2} - 1\right)$$

- We use $a = 3$ as suggested by Liang *et al.* (2008, *JASA*).
- $\frac{\delta}{1+\delta}$ has an interpretation similar to a shrinkage parameter since it accounts for the proportion of information (in data-points) coming from the actual data when $n = n^*$ — in the general case this will be given by $n/(n + n^*/\delta)$.
- Another alternative option would be a hyper- δ/n prior of the form

$$\pi(\delta) = \frac{a-2}{2n} \left(1 + \frac{\delta}{n}\right)^{-a/2}.$$

Additional MCMC step for δ

Step 6: Sample of δ from the full conditional posterior (given the current values of β_γ , β_0 , \mathbf{y}^* and γ).

(a) Propose δ' from $q(\delta'|\delta) = \text{Gamma}(\delta, 1)$.

(b) Compute the Laplace approximations $\hat{m}_\gamma^N(\mathbf{y}^*|\delta)$ and $\hat{m}_\gamma^N(\mathbf{y}^*|\delta')$.

(c) Accept the proposed move with probability $\alpha_\delta = \min\{1, A_\delta\}$, where A_δ is given by

$$A_\delta = \left\{ f_\gamma(\mathbf{y}^*|\beta_\gamma) f_0(\mathbf{y}^*|\beta_0) \right\}^{\Delta\delta} \times \frac{\pi(\delta')}{\pi(\delta)} \times \frac{\hat{m}_\gamma^N(\mathbf{y}^*|\delta)}{\hat{m}_\gamma^N(\mathbf{y}^*|\delta')} \times \frac{q(\delta|\delta')}{q(\delta'\delta)}.$$

where $\Delta\delta = 1/\delta' - 1/\delta$

4 Illustrative examples

- A real life example
- A Poisson simulated study

Illustrative example 2: Pima Indians dataset

- Pima Indians diabetes data set (Ripley, 1996).
- $n = 532$ binary responses on diabetes presence (present=1, not present=0) according to the WHO criteria for signs of diabetes.
- $p = 7$ potential covariates which are listed in Table 1 (see next slide).
- The data also used by Holmes and Held (2006, *Bayesian Analysis*) and Bové and Held (2011, *Bayesian Analysis*).
- Beta-binomial prior on model space.

Covariate	Description
X_1	Number of pregnancies
X_2	Plasma glucose concentration (mg/dl)
X_3	Diastolic blood pressure (mm Hg)
X_4	Triceps skin fold thickness (mm)
X_5	Body mass index (kg/m ²)
X_6	Diabetes pedigree function
X_7	Age

Table 1: Potential predictors in the Pima Indians diabetes data set.

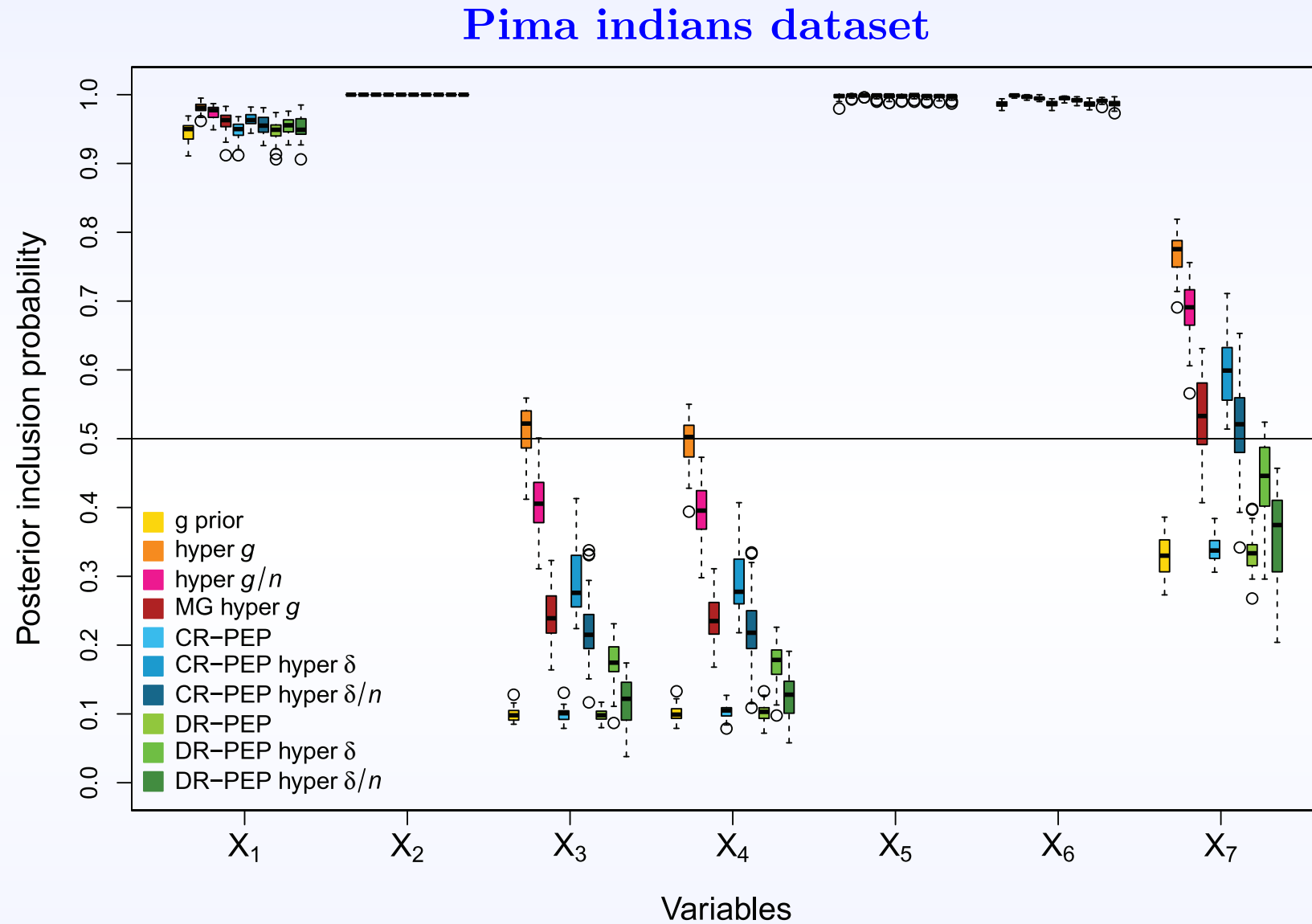


Figure 6: Boxplots of batched estimates of the posterior inclusion probabilities (40 batches of size 1000).

Pima indians dataset

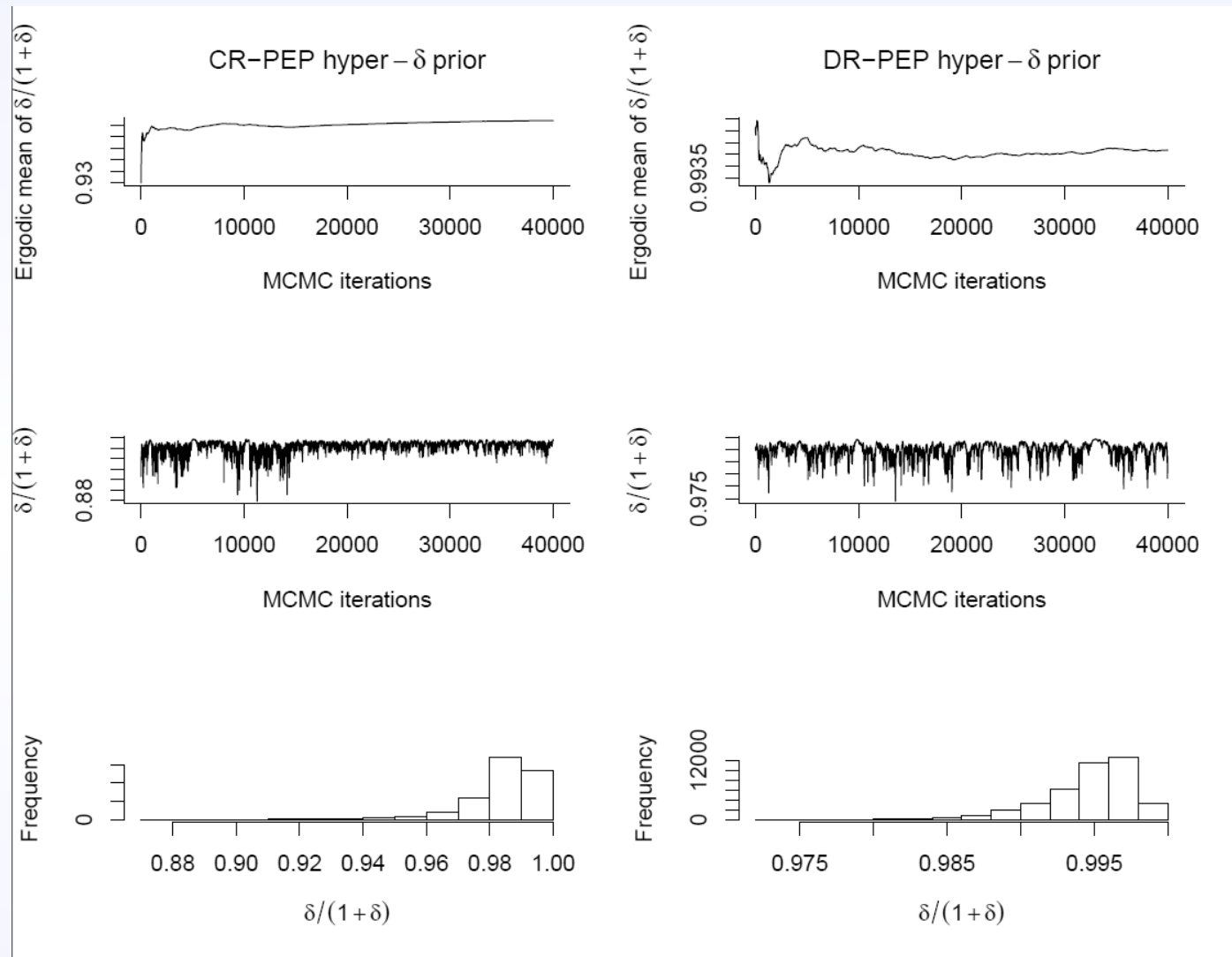


Figure 7: MCMC plots for the shrinkage factor $\delta/(1 + \delta)$ for the CR-PEP and DR-PEP hyper- δ priors (40000 iterations).

Illustrative example 3: Poisson Simulated data

- Also presented in Chen et al. (2008) and Li and Clyde (2015).
- $n = 100$, $p = 3$ predictors. Each simulation is repeated 100 times.
- Each predictor is drawn from a standard normal distribution with pairwise correlation given by

$$\text{corr}(X_i, X_j) = r^{|i-j|}, \quad 1 \leq i < j \leq p.$$

with (i) independent predictors ($r = 0$) and (ii) correlated predictors ($r = 0.75$).

Scenario	Poisson ($n = 100$)			
	β_0	β_1	β_2	β_3
null	-0.3	0	0	0
sparse	-0.3	0.3	0	0
medium	-0.3	0.3	0.2	0
full	-0.3	0.3	0.2	-0.15

Table 2: Four simulation scenarios for Poisson regression assuming independent and correlated predictors.

Prior	Null		Sparse		Medium		Full	
	0	0.75	0	0.75	0	0.75	0	0.75
g -prior	87	93	74	36	29	0	5	0
hyper g -prior	59	71	72	41	45	3	21	2
hyper g/n -prior	81	83	72	42	38	1	13	1
MG hyper g -prior*	84	90	72	37	32	0	10	0
CR PEP	88	95	76	35	27	0	5	0
CR PEP hyper- δ	71	75	68	44	44	4	18	3
CR PEP hyper- δ/n	83	91	80	40	30	0	11	0
DR PEP	90	95	73	32	28	0	5	0
DR PEP hyper- δ	91	97	68	30	25	0	4	0
DR PEP hyper- δ/n	94	95	69	28	20	0	3	0

Table 3: Number of times that the MAP model corresponds to the true model for 100 simulated datasets; column-wise largest value is in red.

Comments on the rates of identifying the true model

- i) Variable selection methods using PEP priors perform well; 6 out of the 8 best MAP success patterns are observed in one of the PEP priors.
- ii) Variable selection methods using PEP priors support more parsimonious models than the competing methods.
- iii) For the **null** and the **sparse** scenarios, PEP priors perform overall better than the competing methods.
- iv) For the **medium** model scenario, the PEP priors perform more or less equally well to the other methods.
- v) When the true model is the **full** model
 - All methods generally fail in the correlated scenario
 - The CR-PEP hyper- δ and δ/n priors are performing generally well in comparison to the competing methods.
 - The rest of the PEP priors have lower MAP success rates than the competing methods using hyper-g priors.

Current directions of research

- We are working to extend the consistency results for the GLMs setup
- **Main direction:** To extend the methodology in *large p , small n* problems.
- Use double exponential as baseline prior
 - g -prior type of behaviour when d_γ is smaller than n
 - LASSO type of shrinkage and behaviour when d_γ is bigger than n or we have extreme collinearities
- What about computation?
EMVS (Rockova and George, 2014, *JASA*) or other fast alternatives should be explored.

Concluding remarks

- We have extended PEP-variable selection for GLMs
- Main problems
 - Definition of the power-likelihood - we have presented two alternatives
 - Computation - we have used an augmented Gibbs variable selection sampler
- CR-PEP and DR-PEP are more parsimonious than g-priors with similar properties.
- Work must be done to prove consistency in the general setup and extend methodology for *large p, small n* problems.



Appendix A: Detailed description of the MCMC algorithm

Step 1: For $j = 1, \dots, p$, update $\gamma_j \sim \text{Bernoulli} \left(\frac{O_j}{1+O_j} \right)$, with

$$O_j = \frac{f_{\gamma_j^1}(\mathbf{y} | \boldsymbol{\beta}_{\gamma_j^1})}{f_{\gamma_j^0}(\mathbf{y} | \boldsymbol{\beta}_{\gamma_j^0})} \times \left[\frac{f_{\gamma_j^1}(\mathbf{y}^* | \boldsymbol{\beta}_{\gamma_j^1})}{f_{\gamma_j^0}(\mathbf{y}^* | \boldsymbol{\beta}_{\gamma_j^0})} \right]^{1/\delta} \times \frac{\pi_{\gamma_j^1}^{\text{N}}(\boldsymbol{\beta}_{\gamma_j^1})}{\pi_{\gamma_j^0}^{\text{N}}(\boldsymbol{\beta}_{\gamma_j^0})} \times \frac{\pi_{\gamma_j^1}(\boldsymbol{\beta}_{\setminus \gamma_j^1})}{\pi_{\gamma_j^0}(\boldsymbol{\beta}_{\setminus \gamma_j^0})} \times \frac{\widehat{m}_{\gamma_j^0}^{\text{N}}(\mathbf{y}^* | \delta)}{\widehat{m}_{\gamma_j^1}^{\text{N}}(\mathbf{y}^* | \delta)} \times \frac{\pi(\gamma_j^1)}{\pi(\gamma_j^0)}.$$

where

- $\gamma_j^1 = (\gamma_j = 1, \gamma_{\setminus j})$
- $\gamma_j^0 = (\gamma_j = 0, \gamma_{\setminus j})$
- $\gamma_{\setminus j}$ is γ without element j

Step 2: Update β_γ from the full conditional posterior (given the current values of γ and \mathbf{y}^*) using a Metropolis step.

- (a) Propose β'_γ from the proposal distribution $q(\beta_\gamma) = \text{N}_{d_\gamma}(\tilde{\beta}_\gamma, \tilde{\Sigma}_{\beta_\gamma})$;
- $\tilde{\beta}_\gamma$ is the MLE with response $(\mathbf{y}, \mathbf{y}^*)$ and weights $\mathbf{w} = (\mathbf{1}_n, \delta^{-1}\mathbf{1}_{n^*})$
 - $\tilde{\Sigma}_{\beta_\gamma}$ the corresponding estimated variance-covariance matrix of $\tilde{\beta}_\gamma$.
- (b) Accept the proposed values with probability $\alpha_{\beta_\gamma} = \min\{1, A_{\beta_\gamma}\}$; where A_{β_γ} is given by

$$A_{\beta_\gamma} = \frac{f_\gamma(\mathbf{y}|\beta'_\gamma)}{f_\gamma(\mathbf{y}|\beta_\gamma)} \times \left[\frac{f_\gamma(\mathbf{y}^*|\beta'_\gamma)}{f_\gamma(\mathbf{y}^*|\beta_\gamma)} \right]^{1/\delta} \times \frac{\pi_\gamma^{\text{N}}(\beta'_\gamma)}{\pi_\gamma^{\text{N}}(\beta_\gamma)} \times \frac{q(\beta_\gamma)}{q(\beta'_\gamma)}.$$

Step 3: Update $\beta_{\setminus\gamma}$ from the pseudo-prior $\pi_\gamma(\beta_{\setminus\gamma}) = \text{N}_{d_{\setminus\gamma}}(\hat{\beta}_{\setminus\gamma}, \mathbf{I}_{d_{\setminus\gamma}} \hat{\sigma}_{\beta_{\setminus\gamma}}^2)$.

Step 4: Sample β_0 from the full conditional posterior (given the current values of \mathbf{y}^*) using a Metropolis step.

- (a) Propose β'_0 from $q(\beta_0) = \text{N}(\hat{\beta}_0^*, \delta \hat{\sigma}_{\hat{\beta}_0^*}^2)$;
- $\hat{\beta}_0^*$ is the MLE with response \mathbf{y}^*

– $\widehat{\sigma}_{\widehat{\beta}_0^*}$ is corresponding standard error of $\widehat{\beta}_0^*$.

(b) Accept the proposed move with probability

$$\alpha_{\beta_0} = \min \left\{ 1, \left[\frac{f_0(\mathbf{y}^* | \beta'_0)}{f_0(\mathbf{y}^* | \beta_0)} \right]^{1/\delta} \times \frac{\pi_0^N(\beta'_0)}{\pi_0^N(\beta_0)} \times \frac{q(\beta_0)}{q(\beta'_0)} \right\}.$$

Step 5: Sample \mathbf{y}^* from the full conditional posterior (given β_γ , β_0 and γ) using a Metropolis step.

(a) Propose $\mathbf{y}^{* \prime}$ from $q(\mathbf{y}^*)$; see Slide 25 for details.

(b) Compute $\widehat{m}_\gamma^N(\mathbf{y}^* | \delta)$ and $\widehat{m}_\gamma^N(\mathbf{y}^{* \prime} | \delta)$ using Laplace approximation.

(c) Accept the proposed values with probability $\alpha_{\mathbf{y}^*} = \min\{1, A_{\mathbf{y}^*}\}$; where $A_{\mathbf{y}^*}$ is given by

$$A_{\mathbf{y}^*} = \left[\frac{f_\gamma(\mathbf{y}^{* \prime} | \beta_\gamma)}{f_\gamma(\mathbf{y}^* | \beta_\gamma)} \times \frac{f_0(\mathbf{y}^{* \prime} | \beta_0)}{f_0(\mathbf{y}^* | \beta_0)} \right]^{1/\delta} \times \frac{\widehat{m}_\gamma^N(\mathbf{y}^* | \delta)}{\widehat{m}_\gamma^N(\mathbf{y}^{* \prime} | \delta)} \times \frac{q(\mathbf{y}^*)}{q(\mathbf{y}^{* \prime})}.$$

Illustrative example 1: Simulated Binomial data (continued)

The following plots also include comparisons with hyper-delta PEP priors for example 1.

Scenario 1 (Null): True model = null

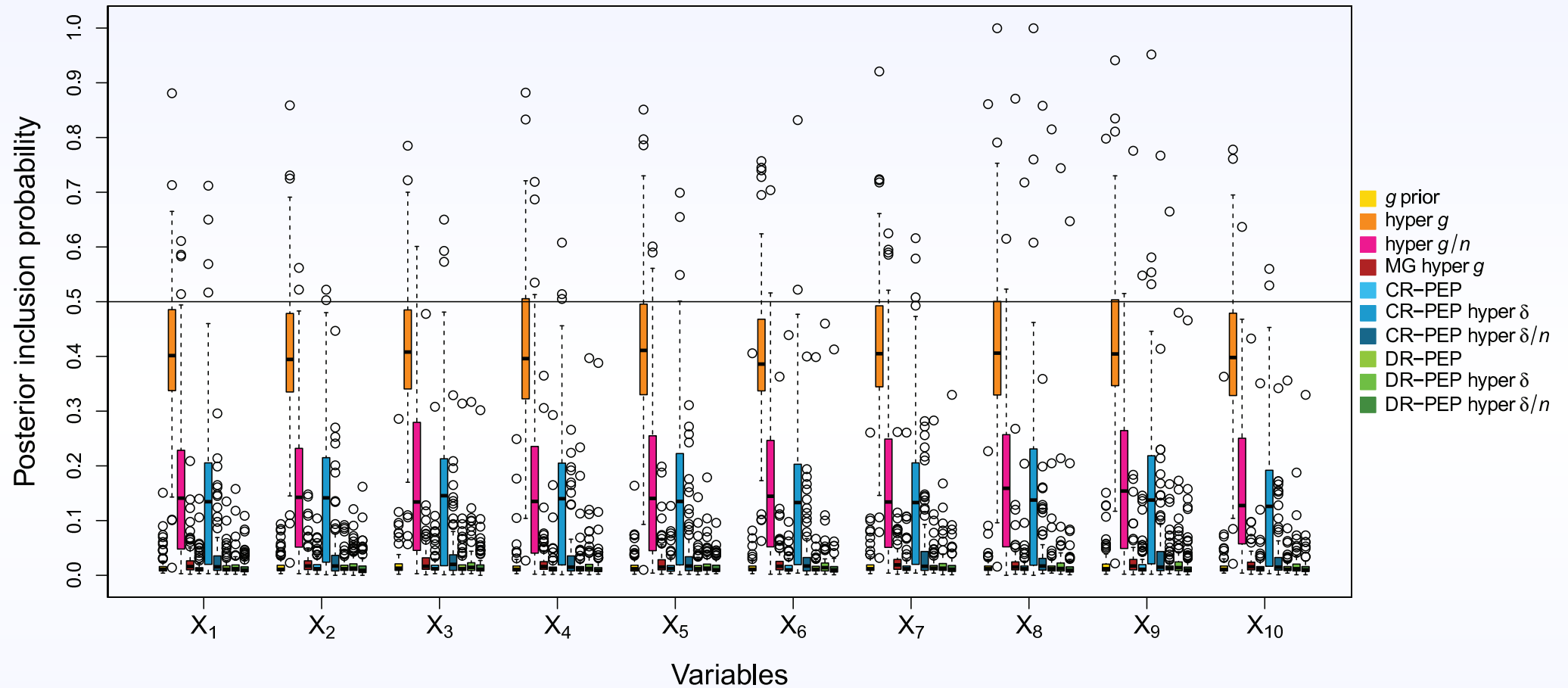


Figure 8: Posterior inclusion probabilities from 100 samples.

MG hyper-g prior: Maruyama & George (2011, Annals of Statistics) prior.

Scenario 2 (Sparse): True model = $0.1 - 0.9X_3 + 1.2X_7 + 0.4X_{10}$

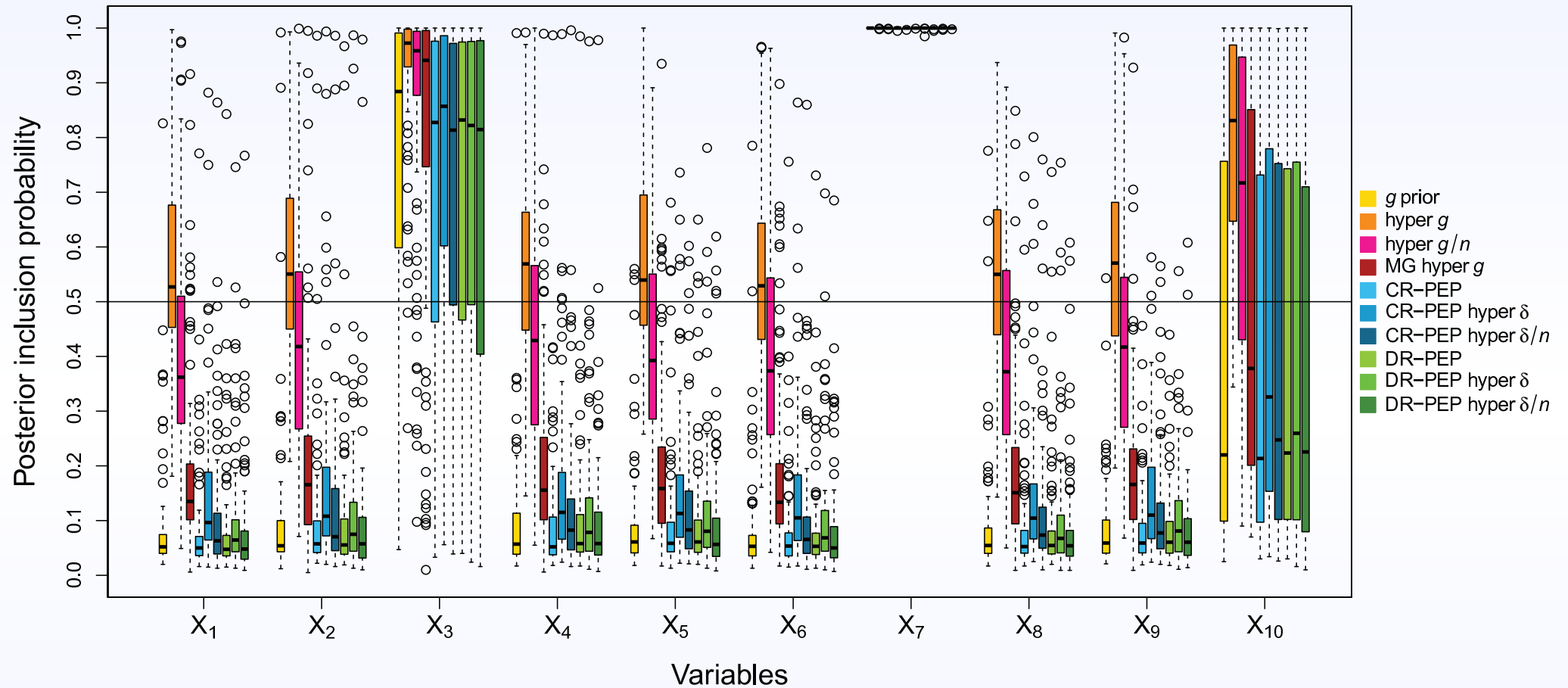


Figure 9: Posterior inclusion probabilities from 100 samples.

MG hyper-g prior: Maruyama & George (2011, Annals of Statistics) prior.

Scenario 3 (Dense)

True model = $0.1 + 0.6X_1 - 0.9X_3 + X_5 + 0.9X_6 + 1.2X_7 - 1.2X_8 - 0.5X_9$

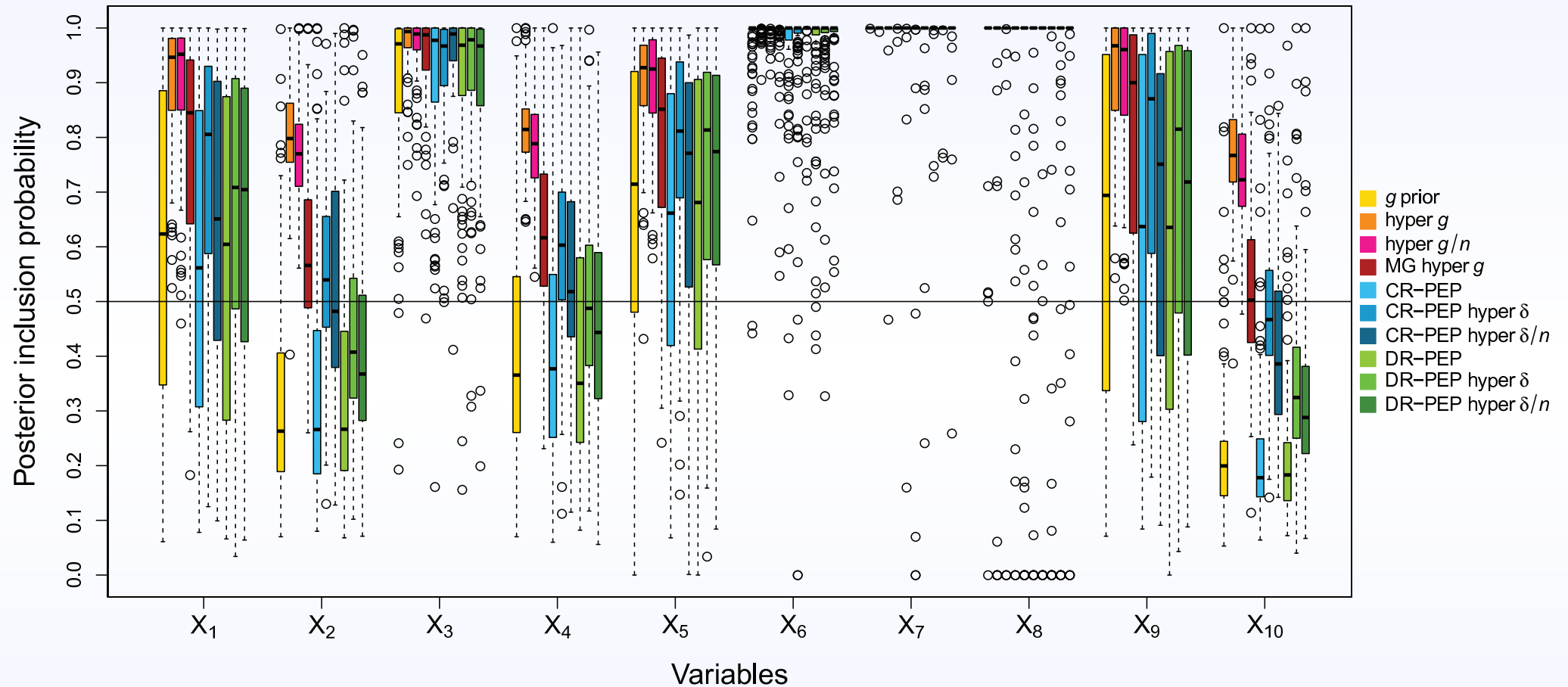


Figure 10: Posterior inclusion probabilities from 100 samples.

MG hyper-g prior: Maruyama & George (2011, Annals of Statistics) prior.

Illustrative example 3: Poisson Simulated data (continued)

Here, we will also find plots of posterior inclusion probabilities for each scenario of simulated example 3.

Independent Covariates

Poisson regression simulation: independent predictors

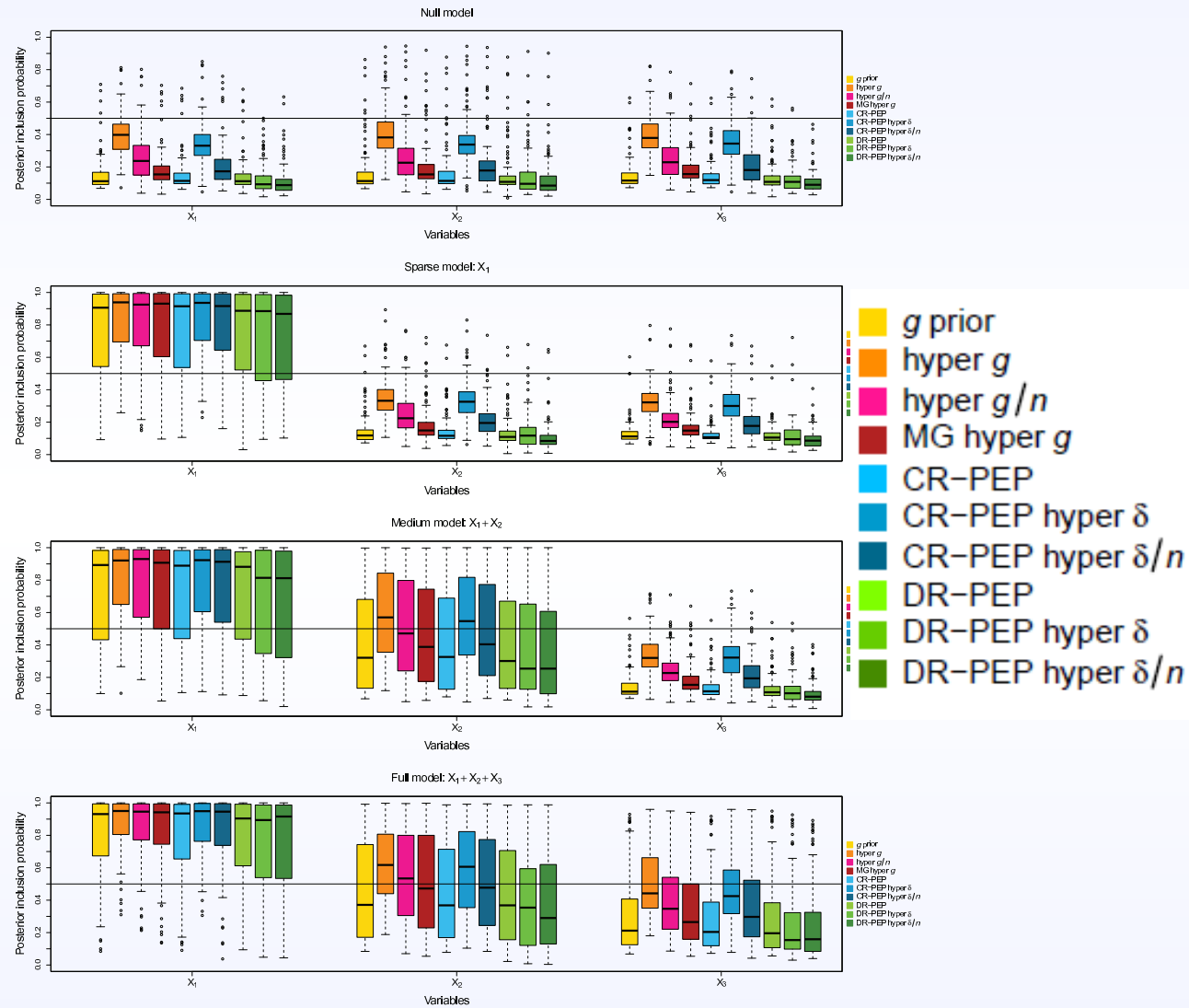


Figure 11: Posterior inclusion probabilities from 100 samples.

Correlated Covariates

Poisson regression simulation: correlated predictors

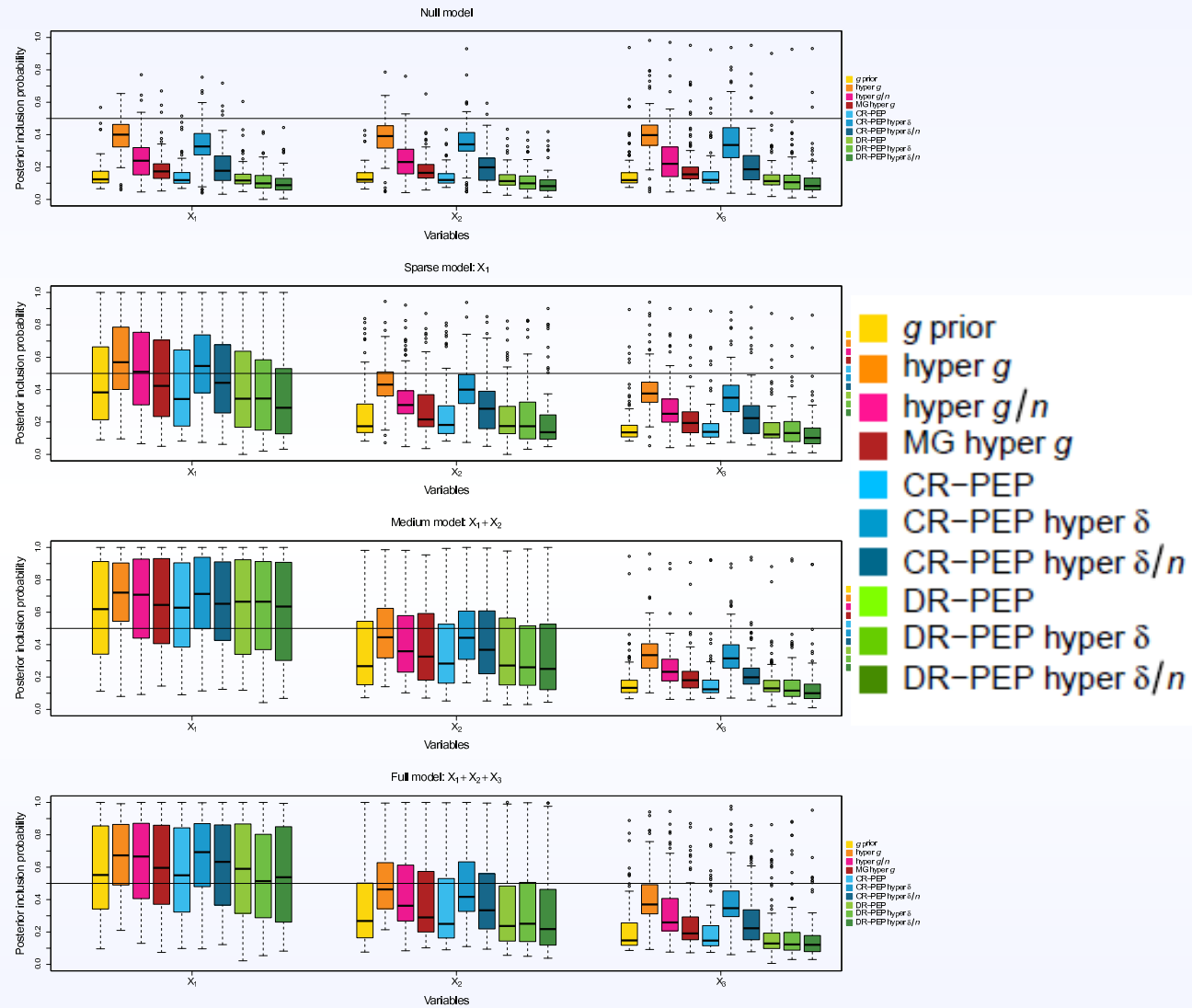


Figure 12: Posterior inclusion probabilities from 100 samples.