

# On the Hyperprior Choice for the Global Shrinkage Parameter in the Horseshoe Prior

Juho Piironen Aki Vehtari

Helsinki Institute for Information Technology, HIIT  
Department of Computer Science, Aalto University, Finland

## INTRODUCTION

- ▶ The horseshoe prior [1] has proven to be a noteworthy choice for **sparse** Bayesian estimation, being a computationally convenient alternative to the **spike-and-slab** prior.
- ▶ The level of sparsity is determined by the **global shrinkage hyperparameter**.
- ▶ However, we demonstrate that the results can be sensitive to the **hyperprior choice** for this parameter.
- ▶ We show how one can specify this hyperprior based on the prior beliefs about the **number of nonzero parameters** in the model.
- ▶ We show that one can **improve the parameter estimation and predictive accuracy** by transforming even a crude prior guess about the sparsity into the model using our framework.

## HORSESHOE PRIOR

- ▶ Consider the standard linear regression model

$$y_i = \beta^T \mathbf{x}_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n,$$

where  $\mathbf{x}$  is the  $D$ -dimensional vector of predictors,  $\beta$  denotes the corresponding coefficients and  $\sigma^2$  is the noise variance.

- ▶ The **horseshoe prior** for the regression coefficients  $\beta = (\beta_1, \dots, \beta_D)$  is given by

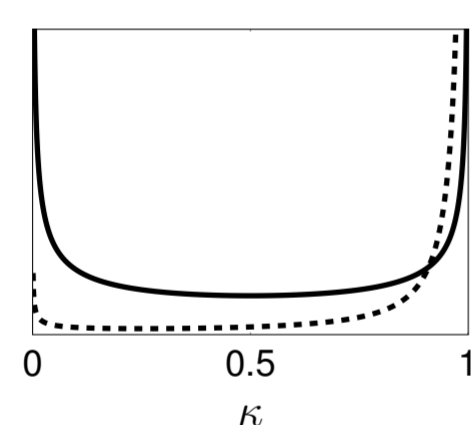
$$\begin{aligned} \beta_j | \lambda_j, \tau &\sim N(0, \lambda_j^2 \tau^2), \\ \lambda_j &\sim C^+(0, 1), \quad j = 1, \dots, D. \end{aligned} \quad (1)$$

- ▶ Given the hyperparameters  $\lambda_j$  and  $\tau$ , and assuming uncorrelated predictors (with zero mean and unit variance), the **posterior mean** satisfies approximately

$$\tilde{\beta}_j = (1 - \kappa_j) \hat{\beta}_j, \quad \kappa_j = \frac{1}{1 + n\sigma^{-2}\tau^2\lambda_j^2}. \quad (2)$$

- ▶ Here  $\hat{\beta}$  is the maximum likelihood solution and  $\kappa_j$  the **shrinkage factor**.

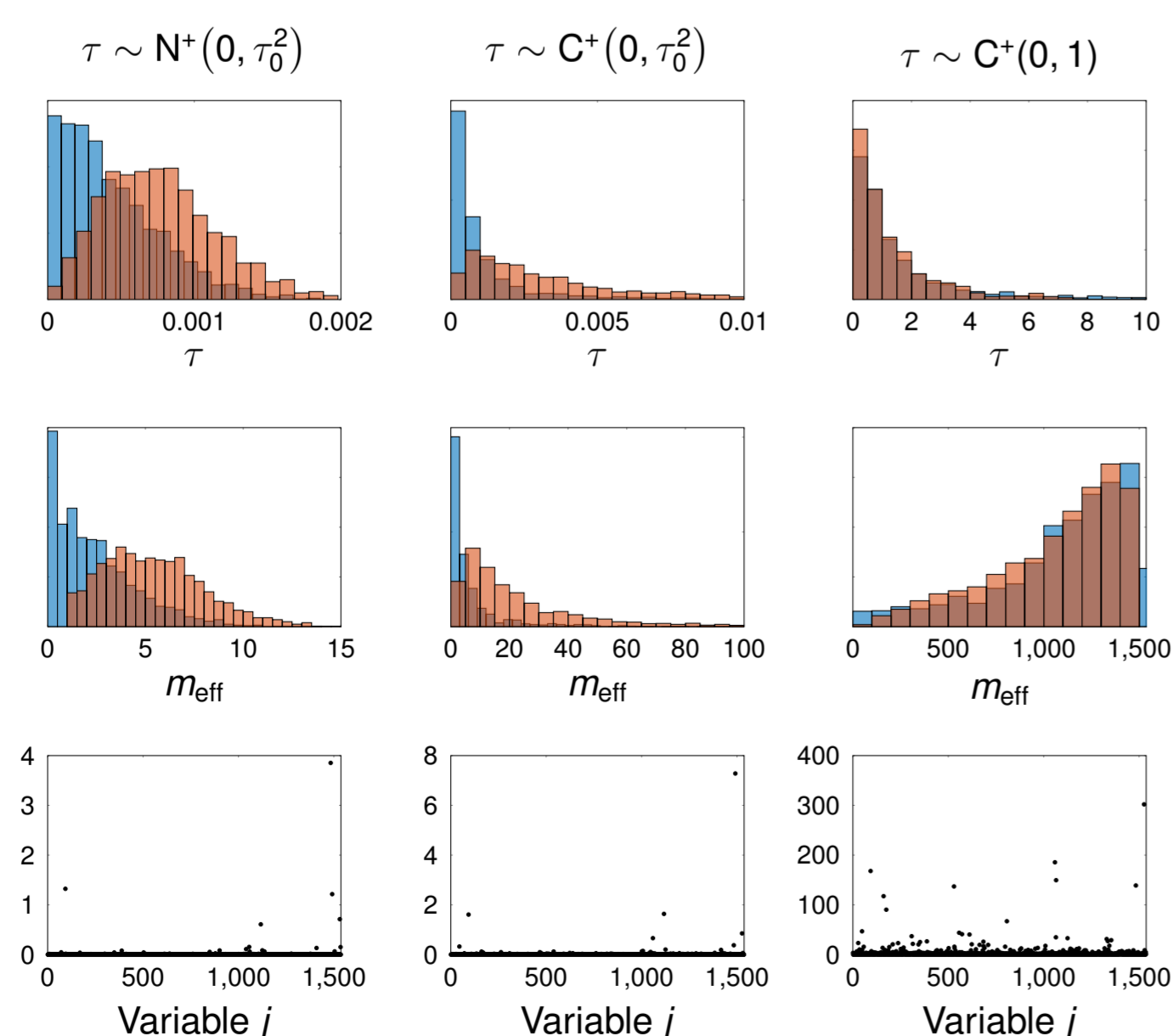
- ▶ The prior for each shrinkage factor  $\kappa_j$  looks like a horseshoe:



Density for the shrinkage factor  $\kappa_j$  (Eq. (2)) for the horseshoe prior (Eq. (1)) when  $n\sigma^{-2}\tau^2 = 1$  (solid) and when  $n\sigma^{-2}\tau^2 = 0.1$  (dashed).

- ▶ Intuition: we expect both relevant ( $\tilde{\beta}_j \approx \hat{\beta}_j$ ) and irrelevant ( $\tilde{\beta}_j \approx 0$ ) variables; which one is favored, depends on the **global shrinkage**  $\tau$ .

## EFFECT ON PARAMETER ESTIMATES



**Ovarian cancer dataset:** Histograms of **prior** and **posterior** samples for  $\tau$  (top row) and  $m_{\text{eff}}$  (middle row), and absolute values of the posterior means for the regression coefficients  $|\tilde{\beta}_j|$  (bottom row) imposed by three different prior choices for  $\tau$ .  $\tau_0$  corresponds to a prior guess  $\rho_0 = 3$  relevant variables (Eq. (4)).

## THE GLOBAL HYPERPARAMETER

- ▶ We define the **effective number of nonzero coefficients** as

$$m_{\text{eff}} = \sum_{j=1}^D (1 - \kappa_j). \quad (3)$$

- ▶ The prior **mean** and **variance** for  $m_{\text{eff}}$  can be derived analytically

$$\begin{aligned} E[m_{\text{eff}} | \tau, \sigma] &= \frac{\tau\sigma^{-1}\sqrt{n}}{1 + \tau\sigma^{-1}\sqrt{n}} D, \\ \text{Var}[m_{\text{eff}} | \tau, \sigma] &= \frac{\tau\sigma^{-1}\sqrt{n}}{2(1 + \tau\sigma^{-1}\sqrt{n})^2} D. \end{aligned}$$

- ▶ Thus, if our prior guess for the number of relevant variables is  $\rho_0$ , it is reasonable to choose the prior so that  $E[m_{\text{eff}} | \tau, \sigma] = \rho_0$ , which yields for  $\tau$

$$\tau_0 = \frac{\rho_0}{D - \rho_0} \frac{\sigma}{\sqrt{n}}. \quad (4)$$

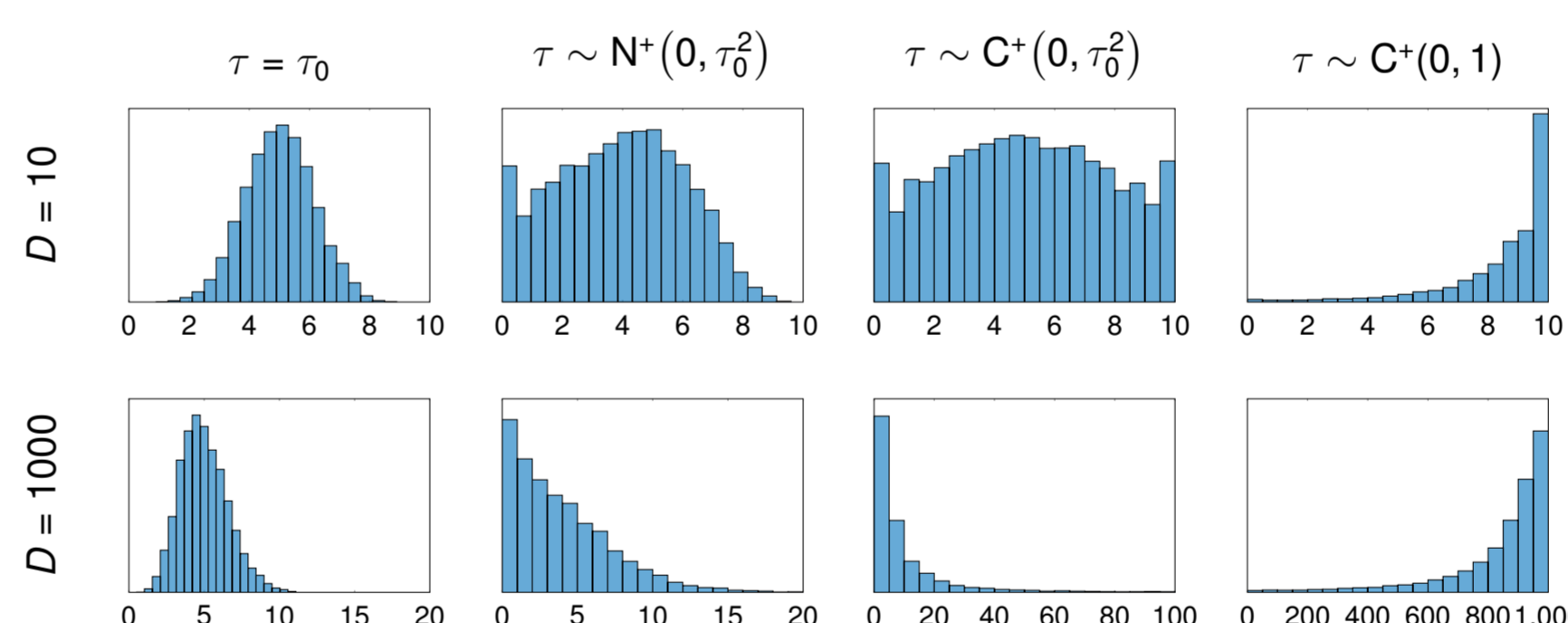
- ▶ This equation captures the relationship between the **global shrinkage parameter** and the **prior assumptions about the sparsity**, and indicates where  $p(\tau)$  should have most of its mass.

- ▶ It is insightful to visualize the prior imposed on  $m_{\text{eff}}$  for different prior choices for  $\tau$  by drawing samples for  $m_{\text{eff}}$  (see the figure below).

- ▶ Theoretical result: assuming a true  $\beta_*$  exists, if our prior guess is  $\rho_0 = p_*$  (the true number of relevant variables), then  $\tau_0$  is asymptotically **the optimal choice** in terms of posterior contraction rates and mean squared error in comparison to the true  $\beta^*$  [2, 3].

- ▶ The result (4) can be generalized also to non-Gaussian observation models by deriving appropriate plug-in values for  $\sigma$ , for instance  $\sigma = 2$  for the logistic regression [3].

## ILLUSTRATION OF THE PRIOR CHOICE



**Various priors for  $m_{\text{eff}}$ :** Histograms of prior draws for  $m_{\text{eff}}$  (Eq. (3)) imposed by different prior choices for  $\tau$ .  $D$  denotes the total number of variables, and  $\tau_0$  is computed from formula (4) assuming  $n = 100$  observations with  $\sigma = 1$  and  $\rho_0 = 5$  as the prior guess for the number of relevant variables. Notice how the “uninformative”  $\tau \sim C^+(0, 1)$  results in a rather dubious prior for  $m_{\text{eff}}$ .

## CONCLUSIONS

- ▶ We have shown how to specify the hyperprior for the **global shrinkage parameter** in the horseshoe prior based on our prior beliefs about the **number nonzero parameters** in the model.
- ▶ Setting up the prior for  $\tau$  based on the prior beliefs regarding the sparsity **improves the results** even when the prior knowledge is rough.
- ▶ The presented framework could also be generalized to **other shrinkage priors** than the horseshoe.

## IMPLEMENTATION

- ▶ The horseshoe prior is implemented in the R-package *rstanarm* (<https://github.com/stan-dev/rstanarm>).
- ▶ A demo about model fitting and the subsequent projective variable selection using our R-package *projpred* (<https://github.com/stan-dev/projpred>) can be found in the vignette <https://users.aalto.fi/~jtpiiron/projpred/quickstart.html>.

## REFERENCES

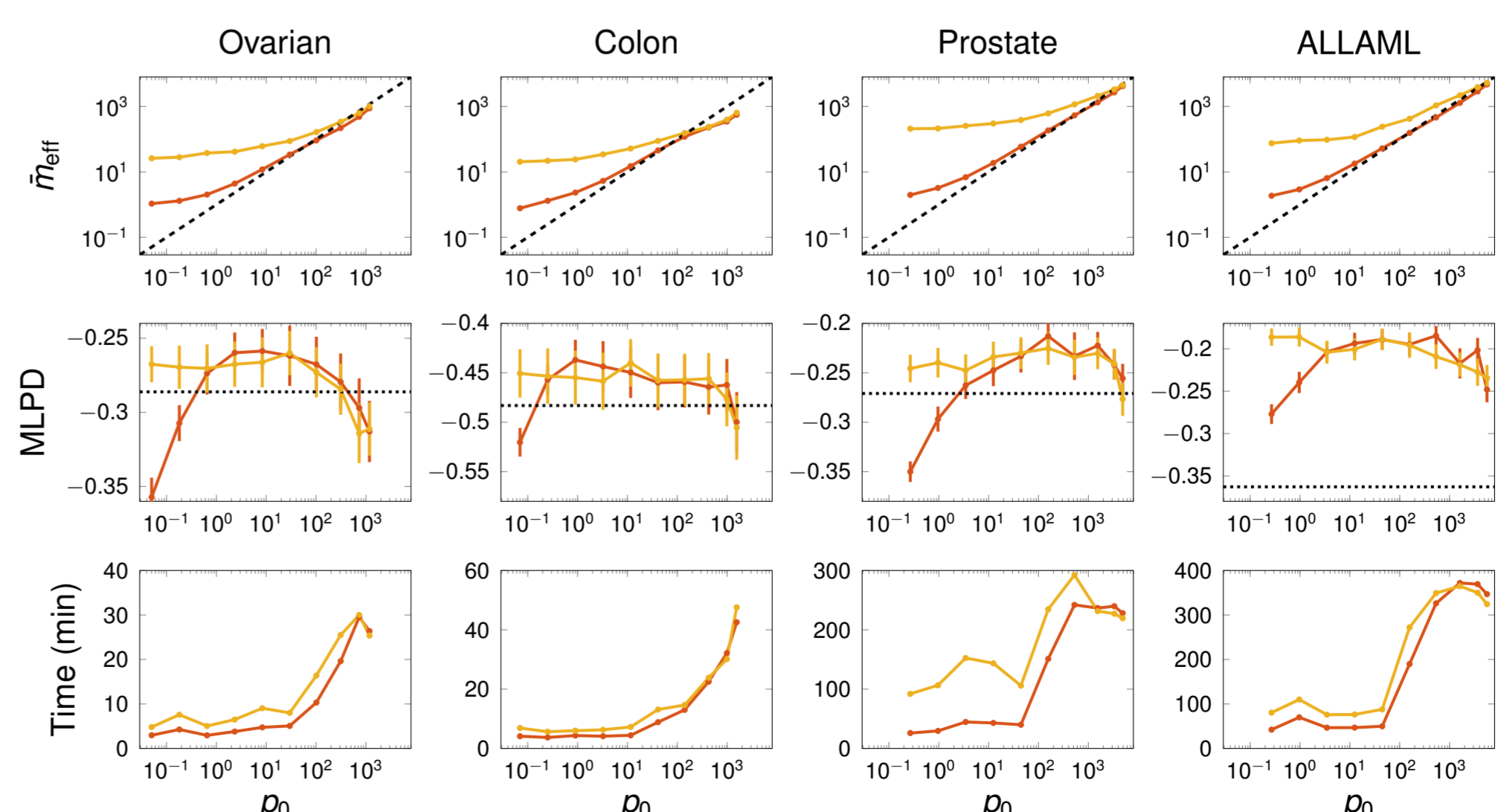
- [1] Carlos M. Carvalho, Nicholas G. Polson, and James G. Scott. Handling sparsity via the horseshoe. In *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics*. JMLR W&CP, 2009. 1
- [2] S. L. van der Pas, B. J. K. Kleijn, and A. W. van der Vaart. The horseshoe estimator: posterior concentration around nearly black vectors. *Electronic Journal of Statistics*, 8(2):2585–2618, 2014. 1
- [3] Juho Piironen and Aki Vehtari. On the hyperprior choice for the global shrinkage parameter in the horseshoe prior. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*. JMLR W&CP, 2017. 1

## DATASETS

| Dataset  | $n$ | $D$  |
|----------|-----|------|
| Ovarian  | 54  | 1536 |
| Colon    | 62  | 2000 |
| Prostate | 102 | 5966 |
| ALLAML   | 72  | 7129 |

Summary of the microarray cancer datasets (binary classification) used for the real world illustrations.

## EFFECT ON PREDICTIVE ACCURACY AND COMPUTATIONAL EFFICIENCY



**Microarray classification datasets:** Posterior mean for  $m_{\text{eff}}$ , mean log predictive density (MLPD) on test data ( $\pm$  one standard error), and computation time for two priors for the global hyperparameter:  $\tau \sim N^+(0, \tau_0^2)$  (red), and  $\tau \sim C^+(0, \tau_0^2)$  (yellow), where  $\tau_0$  is computed from (4) varying  $\rho_0$  (horizontal axis). For each curve, the largest  $\rho_0$  corresponds to  $\tau_0 = 1$ . For comparison, the dotted line in the middle row plots denotes the MLPD for LASSO.