

# Variable Selection with ABC Bayesian Forests

Yi Liu

*Department of Statistics, University of Chicago*

Veronika Ročková

*Booth School of Business, University of Chicago*

Yuexi Wang

*Booth School of Business, University of Chicago*

**Summary.** Few problems in statistics are as perplexing as variable selection in the presence of very many redundant covariates. The variable selection problem is most familiar in parametric environments such as the linear model or additive variants thereof. In this work, we abandon the linear model framework, which can be quite detrimental when the covariates impact the outcome in a non-linear way, and turn to tree-based methods for variable selection. Such variable screening is traditionally done by pruning down large trees or by ranking variables based on some importance measure. Despite heavily used in practice, these ad-hoc selection rules are not yet well understood from a theoretical point of view. In this work, we devise a Bayesian tree-based probabilistic method and show that it is consistent for variable selection when the regression surface is a smooth mix of  $p > n$  covariates. These results are the first model selection consistency results for Bayesian forest priors. Probabilistic assessment of variable importance is made feasible by a spike-and-slab wrapper around sum-of-trees priors. Sampling from posterior distributions over trees is inherently very difficult. As an alternative to MCMC, we propose ABC Bayesian Forests, a new ABC sampling method based on data-splitting that achieves higher ABC acceptance rate. We show that the method is robust and successful at finding variables with high marginal inclusion probabilities. Our ABC algorithm provides a new avenue towards approximating the median probability model in non-parametric setups where the marginal likelihood is intractable.

*Keywords:* Approximate Bayesian Computation, BART, Consistency, Spike-and-Slab, Variable Selection

## 1. Perspectives on Non-parametric Variable Selection

In its simplest form, variable selection is most often carried out in the context of linear regression ([Tibshirani, 1996](#); [George and McCulloch, 1993](#); [Fan and Li, 2001](#)). However, confinement to linear parametric forms can be quite detrimental for variable importance screening, when the covariates impact the outcome in a non-linear way ([Turlach, 2004](#)). Rather than first selecting a parametric model to filter out variables, another strategy is to first select variables and then build a model. Adopting this reversed point of view, we focus on developing methodology for the so called “model-free” variable selection ([Chipman et al., 2001](#)).

There is a long strand of literature on the fundamental problem of non-parametric variable selection. One line of research focuses on capturing non-linearities and interactions with basis expansions and performing grouped shrinkage/selection on sets of coefficients (Scheipl, 2011; Ravikumar et al., 2009; Lin and Zhang, 2006; Radchenko and James, 2010). Lafferty and Wasserman (2008) propose the RODEO method for sparse non-parametric function estimation through regularization of the derivative expectation operator and provide a consistency result for the selection of the optimal bandwidth. Candès et al. (2018) propose a model-free knock-off procedure, controlling FDR in settings when the conditional distribution of the response is arbitrary. In the Bayesian literature, Savitsky et al. (2011) deploy spike-and-slab priors on covariance parameters of Gaussian processes to erase variables. In this work, we focus on other non-parametric regression techniques, namely trees/forests which have been ubiquitous throughout machine learning and statistics (Breiman, 2001; Chipman et al., 2010). The question we wish to address is whether one can leverage the flexibility of regression trees for effective (consistent) variable importance screening.

While trees are routinely deployed for data exploration, prediction and causal inference (Hill, 2011; Taddy et al., 2011a; Gramacy and Lee, 2008), they have also been used for dimension reduction and variable selection. This is traditionally done by pruning out variables or by ranking them based on some importance measure. The notion of variable importance was originally proposed for CART using overall improvement in node impurity involving surrogate predictors (Breiman et al., 1984). In random forests, for example, the importance measure consists of a difference between prediction errors before and after noising the covariate through a permutation in the out-of-bag sample. However, this continuous variable importance measure is on an arbitrary scale, rendering variable selection ultimately ad-hoc. Principled selection of the importance threshold (with theoretical guarantees such as FDR control or model selection consistency) is still an open problem. Simplified variants of importance measures have begun to be understood theoretically for variable selection only very recently (Ishwaran, 2007; Kazemitabar et al., 2017).

Bayesian trees and forests select variables based on probabilistic considerations. The BART procedure (Chipman et al., 2010) can be adapted for variable selection by forcing the number of available splits (trees) to be small, thereby introducing competition between predictors. BART then keeps track of predictor inclusion frequencies and outputs a probabilistic importance measure: an average proportion of all splitting rules inside a tree ensemble that split on a given variable, where the average is taken over the MCMC samples. This measure cannot be directly interpreted as the posterior variable inclusion probability in anisotropic regression surfaces, where wigglier directions require more splits. Bleich et al. (2014) consider a permutation framework for obtaining the null distribution of the importance weights. Zhu et al. (2015) implement reinforcement learning for selection of splitting variables during tree construction to encourage splits on fewer more important variables. All these developments point to the fact that regularization is key to enhancing performance of trees/forests in high dimensions. Our approach differs in that we impose regularization from *outside* the tree/forest through a spike-and-slab wrapper.

Spike-and-slab variable selection consistency results have relied on analytical tractabil-

ity (approximation availability) of the marginal likelihood (Narisetty and He, 2014; Johnson and Rossell, 2012; Castillo et al., 2015). Nicely tractable marginal likelihoods are ultimately unavailable in our framework, rendering the majority of the existing theoretical tools inapplicable. For these contexts, Yang and Pati (2017) characterized general conditions for model selection consistency, extending the work of Lember and van der Vaart (2007) to non *iid* setting. Exploiting these developments, we show variable selection consistency of our non-parametric spike-and-slab approach when the regression function is a smooth mix of covariates. Building on Ročková and van der Pas (2017), our paper continues the investigation of missing theoretical properties of Bayesian CART and BART. We show model selection consistency when the smoothness is known as well as joint consistency for both the regularity level *and* active variable set when the smoothness is not known and when  $p > n$ . These results are the first model selection consistency results for Bayesian forest priors.

The absence of a tractable marginal likelihood complicates not only theoretical analysis, but also computation. We turn to Approximate Bayesian Computation (ABC) (Plagnol and Tavaré, 2004; Marin et al., 2012; Csillery et al., 2010) and propose a procedure for model-free variable selection. Our ABC method *does not* require the use of low-dimensional summary statistics and, as such, it *does not* suffer from the known difficulty of ABC model choice (Robert et al., 2011). Our method is based on sample splitting where at each iteration (a) a random subset of data is used to come up with a proposal draw and (b) the rest of the data is used for ABC acceptance. This new data-splitting approach increases ABC effectiveness by increasing its acceptance rate. ABC Bayesian forests relate to the recent line of work on combining machine learning with ABC (Pudlo et al., 2015; Jiang et al., 2017). We propose dynamic plots that describe the evolution of marginal inclusion probabilities as a function of the ABC selection threshold.

The paper is structured as follows. Section 2 introduces the spike-and-slab wrapper around tree priors. Section 3 develops the ABC variable selection algorithm. Section 4 presents model selection consistency results. Section 5 demonstrates the usefulness of the ABC method on simulated data and Section 6 wraps up with a discussion.

### 1.1. Notation

With  $\|\cdot\|_n$  we denote the empirical  $L^2$  norm. The class of functions  $f(\mathbf{x}) : [0, 1]^p \rightarrow \mathbb{R}$  such that  $f(\cdot)$  is constant in all directions excluding  $\mathcal{S}_0 \subseteq \{1, \dots, p\}$  is denoted with  $\mathcal{C}(\mathcal{S}_0)$ . With  $\mathcal{H}_p^\alpha$ , we denote  $\alpha$ -Hölder continuous functions with a smoothness coefficient  $\alpha$ .  $a \lesssim b$  denotes  $a$  is less or equal to  $b$ , up to a multiplicative positive constant, and  $a \asymp b$  denotes  $a \lesssim b$  and  $b \lesssim a$ . The  $\varepsilon$ -covering number of a set  $\Omega$  for a semimetric  $d$ , denoted by  $N(\varepsilon; \Omega; d)$ , is the minimal number of  $d$ -balls of radius  $\varepsilon$  needed to cover set  $\Omega$ .

## 2. Bayesian Subset Selection with Trees

We will work within the purview of non-parametric regression, where a vector of continuous responses  $\mathbf{Y}^{(n)} = (Y_1, \dots, Y_n)'$  is linked to fixed (rescaled) predictors  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})' \in [0, 1]^p$  for  $1 \leq i \leq n$  through

$$Y_i = f_0(\mathbf{x}_i) + \varepsilon_i \quad \text{with} \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \text{for} \quad 1 \leq i \leq n, \quad (1)$$

where  $f_0(\cdot)$  is the regression mixing function and  $\sigma^2 > 0$  is a scalar. It is often reasonable to expect that only a small subset  $\mathcal{S}_0$  of  $q_0 = |\mathcal{S}_0|$  predictors actually exert influence on  $\mathbf{Y}^{(n)}$  and contribute to the mix. The subset  $\mathcal{S}_0$  is seldom known with certainty and we are faced with the problem of variable selection. Throughout this paper, we assume that the regression surface is smoothly varying ( $\alpha$ -Hölder continuous) along the active directions  $\mathcal{S}_0$  and constant otherwise, i.e. we write  $f_0 \in \mathcal{H}_p^\alpha \cap \mathcal{C}(\mathcal{S}_0)$ .

Unlike linear models that capture the effect of a single covariate with a single coefficient, we permit non-linearities/interactions and capture variable importance with (additive) regression trees. By doing so, we hope to recover non-linear signals that could be otherwise missed by linear variable selection techniques.

As with any other non-parametric regression method, regression trees are vulnerable to the curse of dimensionality, where prediction performance deteriorates dramatically as the number of variables  $p$  increases. If an oracle were to isolate the active covariates  $\mathcal{S}_0$ , the fastest achievable estimation rate would be  $n^{-\alpha/(2\alpha+|\mathcal{S}_0|)}$ . This rate depends only on the intrinsic dimensionality  $q_0 = |\mathcal{S}_0|$ , not the actual dimensionality  $p$  which can be much larger than  $n$ . Recently, [Ročková and van der Pas \(2017\)](#) showed that with *suitable regularization*, the posterior distribution for Bayesian CART and BART actually concentrates at this fast rate (up to a log factor), adapting to the intrinsic dimensionality and smoothness. Later in [Section 4](#), we continue their theoretical investigation and focus on consistent *variable selection*, i.e. estimation of  $\mathcal{S}_0$  rather than  $f_0(\cdot)$ . Spike-and-slab regularization plays a key role in obtaining these theoretical guarantees.

### 2.1. Trees with Spike-and-Slab Regularization

Many applications offer a plethora of predictors and some form of redundancy penalization has to be incurred to cope with the curse of dimensionality. Bayesian regression trees were originally conceived for prediction rather than variable selection. Indeed, original tree implementations of Bayesian CART ([Denison et al., 1998](#); [Chipman et al., 1998](#)) do not seem to penalize inclusion of redundant variables aggressively enough. As noted by [Linero \(2018\)](#), the prior expected number of active variables under the Bayesian CART prior of [Chipman et al. \(1998\)](#) satisfies  $\lim_{p \rightarrow \infty} \mathbb{E}[q] = K - 1$  as  $p \rightarrow \infty$  where  $K$  is the fixed number of bottom leaves. This behavior suggests that (in the limit) the prior forces inclusion of the maximal number of variables while splitting on them only once. This is far from ideal. To alleviate this issue, we deploy the so-called *spike-and-forest priors*, i.e. spike-and-slab wrappers around sum-of-trees priors ([Ročková and van der Pas, 2017](#)). As with the traditional spike-and-slab priors, the specification starts with a prior distribution over the  $2^p$  active variable sets:

$$\mathcal{S} \sim \pi(\mathcal{S}) \quad \text{for each } \mathcal{S} \subseteq \{1, \dots, p\}. \quad (2)$$

We elaborate on the specific choices of  $\pi(\mathcal{S})$  later in [Section 3.2](#) and [Section 4](#).

Given the pool of variables  $\mathcal{S}$ , a regression tree/forest is grown using *only* variables inside  $\mathcal{S}$ . This prevents the trees from using too many variables and thereby from overfitting. Recall that each individual regression tree is characterized by two components: (1) a tree-shaped  $K$ -partition of  $[0, 1]^p$ , denoted with  $\mathcal{T}$ , and (2) bottom node parameters (step heights), denoted with  $\boldsymbol{\beta} \in \mathbb{R}^K$ . Starting with a parent node  $[0, 1]^p$ , each  $K$ -partition is grown by recursively dissecting rectangular cells at chosen internal nodes

along one of the active coordinate axes, all the way down to  $K$  terminal nodes. Each tree-shaped  $K$ -partition  $\mathcal{T} = \{\Omega_k\}_{k=1}^K$  consists of  $K$  partitioning rectangles  $\Omega_k \subset [0, 1]^p$ .

While Bayesian CART approximates  $f_0(\mathbf{x})$  with a single tree mappings  $f_{\mathcal{T}, \beta}(\mathbf{x}) = \sum_{k=1}^K \mathbb{I}(\mathbf{x} \in \Omega_k) \beta_k$ , Bayesian Additive Regression Trees (BART) use an aggregate of  $T$  mappings

$$f_{\mathcal{E}, \mathbf{B}}(\mathbf{x}) = \sum_{t=1}^T f_{\mathcal{T}^t, \beta^t}(\mathbf{x})$$

where  $\mathcal{E} = \{\mathcal{T}^1, \dots, \mathcal{T}^T\}$  is an ensemble of tree partitions and  $\mathbf{B} = [\beta^1, \dots, \beta^T]$  is an ensemble of step coefficients. In a fully Bayesian approach, prior distributions have to be specified over the set of tree structures  $\mathcal{E}$  and over terminal node heights  $\mathbf{B}$ . The spike-and-forest construction can accommodate various tree prior options.

To assign a prior over  $\mathcal{E}$  for a given  $T$ , one possibility is to first pick the number of bottom nodes, independently for each tree, from a prior

$$K^t \sim \pi(K) \quad \text{for } K = 1, \dots, n, \quad (3)$$

such as the Poisson distribution (Denison et al., 1998). Given the vector of tree sizes  $\mathbf{K} = (K^1, \dots, K^T)'$  and a set of covariates  $\mathcal{S}$ , we assign a prior over so-called valid ensembles/forests  $\mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}}$ . We say that a tree ensemble  $\mathcal{E}$  is valid if it consists of trees that have non-empty bottom leaves. One can pick a tree partition ensemble from a uniform prior over *valid* forests  $\mathcal{E} \in \mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}}$ , i.e.

$$\pi(\mathcal{E} \mid \mathcal{S}, \mathbf{K}) = \frac{1}{\Delta(\mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}})} \mathbb{I}(\mathcal{E} \in \mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}}), \quad (4)$$

where  $\Delta(\mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}})$  is the number of valid tree ensembles characterized by  $\mathbf{K}$  bottom leaves and split directions  $\mathcal{S}$ . The prior (3) and (4) was deployed in the Bayesian CART implementation of Denison et al. (1998) (with  $T = 1$ ) and it was studied theoretically by Ročková and van der Pas (2017). Another related Bayesian forest prior (implemented in the BART procedure and studied theoretically by Ročková and Saha (2019) consists of an independent product of branching process priors (one for each tree) with decaying split probabilities (Chipman et al., 1998). The implementation is very similar to the one of Denison et al. (1998).

Finally, given the partitions  $\mathcal{T}^t$  of size  $K^t$  for  $1 \leq t \leq T$ , one assigns (independently for each tree) a Gaussian product prior on the step heights

$$\pi(\beta^t \mid K^t) = \prod_{k=1}^{K^t} \phi(\beta_k^t; \sigma_{\beta}^2), \quad (5)$$

where  $\phi(x; \sigma_{\beta}^2)$  denotes a Gaussian density with mean zero and variance  $\sigma_{\beta}^2 = 1/T$  (as suggested by Chipman et al. (2010)). The prior for  $\sigma^2$  can be chosen as inverse chi-squared with hyperparameters chosen based on an estimate of the residual standard deviation of the data (Chipman et al., 2010).

The most crucial component in the spike-and-forest construction, which sets it apart from existing BART implementations, is the active set  $\mathcal{S}$  which serves to mute variables

by restricting the pool of predictors available for splits. The goal is to learn which set  $\mathcal{S}$  is most likely (a posteriori) and/or how likely each variables is to have contributed to  $f_0$ . Unlike related tree-based variable selection criteria, the spike-and-slab envelope makes it possible to perform variable selection directly by evaluating posterior model probabilities  $\Pi(\mathcal{S} | \mathbf{Y}^{(n)})$  or marginal inclusion probabilities  $\Pi(j \in \mathcal{S}_0 | \mathbf{Y}^{(n)})$  for  $1 \leq j \leq p$ . Random forests (Breiman, 2001) also mute variables, but they do so from within the tree by randomly choosing a small subset of variables for each split. The spike-and-slab approach mutes variables externally rather than internally. Bleich et al. (2014) note that when the number of trees is small, the Gibbs sampler for BART can get trapped in local modes which can destabilize the estimation procedure. On the other hand, when the number of trees is large, there are ample opportunities for the noise variables to enter the model without necessarily impacting the model fit, making variable selection very challenging. Our spike-and-slab wrapper is devised to get around this problem.

The problem of variable selection is fundamentally challenged by the sheer size of possible variable subsets. For linear regression, (a) MCMC implementations exist that capitalize on the availability of marginal likelihood (Narisetty and He, 2014; Guan and Stephens, 2011), (b) optimization strategies exist for both continuous (Ročková and George, 2018; Ročková, 2017) and point-mass spike-and slab priors (Carbonetto and Stephens, 2012). These techniques do not directly translate to tree models, for which tractable marginal likelihoods  $\pi(\mathbf{Y}^{(n)} | \mathcal{S})$  are unavailable. To address this computational challenge, we explore ABC techniques as a new promising avenue for non-parametric spike-and-slab methods.

### 3. ABC for Variable Selection

Performing (approximate) posterior inference in complex models is often complicated by the analytical intractability of the marginal likelihood. Approximate Bayesian Computation (ABC) is a simulation-based inference framework that obviates the need to compute the likelihood directly by evaluating the proximity of (sufficient statistics of) observed data and pseudo-data simulated from the likelihood. Simon Tavaré first proposed the ABC algorithm for posterior inference (Tavaré et al., 1997) in the 1990’s and since then it has widely been used in population genetics, systems biology, epidemiology and phylogeography<sup>1</sup>.

Combined with a probabilistic structure over models, marginal likelihoods give rise to posterior model probabilities, a standard tool for Bayesian model choice. When the marginal likelihood is unavailable (our case here), ABC offers a unique computational solution. However, as pointed out by Robert et al. (2011), ABC cannot be trusted for model comparisons when model-wise sufficient summary statistics are not sufficient across models. The ABC approximation to Bayes factors then does not converge to exact Bayes factors, rendering ABC model choice fundamentally untrustworthy. A fresh new perspective to ABC model choice was offered in Pudlo et al. (2015), who rephrase model selection as a classification problem that can be tackled with machine learning tools. Their idea is to treat the ABC reference table (consisting of samples from a prior model distribution and high-dimensional vectors of summary statistics of pseudo-data obtained from the prior predictive distribution) as an actual data set, and to train a random

<sup>1</sup>The study of how human beings migrated throughout the world in the past.



forest classifier that predicts a model label using the summary statistics as predictors. Their goal is to produce a stable model decision based on a classifier rather than on an estimate of posterior model probabilities. Our approach has a similar flavor in the sense that it combines machine learning with ABC, but the concept is fundamentally very different. Here, the fusion of Bayesian forests and ABC is tailored to non-parametric variable selection towards obtaining posterior variable inclusion probabilities. Our model selection approach does not suffer from the difficulty of ABC model choice as we *do not* commit to any summary statistics and use random subsets of observations to generate the ABC reference table.

### 3.1. Naive ABC Implementation

For its practical implementation, our Bayesian variable selection method requires sampling from the analytically intractable posterior distribution over subsets  $\Pi(\mathcal{S} | \mathbf{Y}^{(n)})$  under the *spike-and-forest* prior (4), (3) and (2). Given a single tree partition  $\mathcal{T}$ , the (conditional) marginal likelihood  $\pi(\mathbf{Y}^{(n)} | \mathcal{T}, \mathcal{S})$  is available in closed form, facilitating implementations of Metropolis-Hastings algorithms (Chipman et al., 1998; Denison et al., 1998) (see Section S.3). However, such MCMC schemes can suffer from poor mixing. Taking advantage of the fact that, despite being intractable, one can *simulate from* the marginal likelihood  $\pi(\mathbf{Y}^{(n)} | \mathcal{S})$ , we will explore the potential of ABC as a complementary development to MCMC implementations.

The principle at the core of ABC is to perform approximate posterior inference from a given dataset by simulating from a prior distribution and by comparisons with numerous synthetic datasets. In its standard form, an ABC implementation of model choice creates a reference table, recording a large number of datasets simulated from the model prior and the prior predictive distribution under each model. Here, the table consists of  $M$  pairs  $(\mathcal{S}_m, \mathbf{Y}_m^*)$  of model indices  $\mathcal{S}_m$ , simulated from the prior  $\pi(\mathcal{S})$ , and pseudo-data  $\mathbf{Y}_m^* \in \mathbb{R}^n$ , simulated from the marginal likelihood  $\pi(\mathbf{Y}^{(n)} | \mathcal{S}_m)$ . To generate  $\mathbf{Y}_m^*$  in our setup, one can hierarchically decompose the marginal likelihood

$$\pi(\mathbf{Y}^{(n)} | \mathcal{S}) = \int_{(f_{\mathcal{E}, \mathcal{B}}, \sigma^2)} \pi(\mathbf{Y}^{(n)} | f_{\mathcal{E}, \mathcal{B}}, \sigma^2) d\pi(f_{\mathcal{E}, \mathcal{B}}, \sigma^2 | \mathcal{S}) \quad (6)$$

and first draw  $(f_{\mathcal{E}, \mathcal{B}}^m, \sigma_m^2)$  from the prior  $\pi(f_{\mathcal{E}, \mathcal{B}}, \sigma^2 | \mathcal{S})$  and obtain  $\mathbf{Y}_m^*$  from (1), given  $(f_{\mathcal{E}, \mathcal{B}}^m, \sigma_m^2)$ . ABC sampling is then followed by an ABC rejection step, which extracts pairs  $(\mathcal{S}_m, \mathbf{Y}_m^*)$  such that  $\mathbf{Y}_m^*$  is close enough to the actual observed data. In other words, one trims the reference table by keeping only model indices  $\mathcal{S}_m$  paired with pseudo-observations that are at most  $\epsilon$ -away from the observed data, i.e.  $\|\mathbf{Y}^{obs} - \mathbf{Y}_m^*\|_2 \leq \epsilon$  for some tolerance level  $\epsilon$ . These extracted values comprise an approximate ABC sample from the posterior  $\pi(\mathcal{S} | \mathbf{Y}^{(n)})$ , which should be informative for the relative ordering of the competing models, and thus variable selection (Grelaud et al., 2009). Note that this particular ABC implementation does not require any use of low-dimensional summary statistics, where rejection is based solely on  $\mathbf{Y}^{obs}$ . While theoretically justified, this ABC variant has two main drawbacks.

First, with very many predictors, it will be virtually impossible to sample from all  $2^p$  model combinations at least once, unless the reference table is huge. Consequently, relative frequencies of occurrence of a model  $\mathcal{S}_m$  in the trimmed ABC reference table

may not be a good estimate of the posterior model probability  $\pi(\mathcal{S}_m | \mathbf{Y}^{(n)})$ . While the model with the highest posterior probability  $\pi(\mathcal{S}_m | \mathbf{Y}^{(n)})$  is commonly conceived as the right model choice, it may not be the optimal model for prediction. Indeed, in nested correlated designs and orthogonal designs, it is the median probability model that is predictive optimal (Barbieri and Berger, 2004). The median probability model (MPM) consists of those variables whose *marginal* inclusion probabilities  $\mathbb{P}(j \in \mathcal{S}_0 | \mathbf{Y}^{(n)})$  are at least 0.5. While simulation-based estimates of posterior model probabilities  $\mathbb{P}(\mathcal{S} | \mathbf{Y}^{(n)})$  can be imprecise, we argue (and show) that ABC estimates of marginal inclusion probabilities  $\mathbb{P}(j \in \mathcal{S}_0 | \mathbf{Y}^{(n)})$  are far more robust and stable.

The second difficulty is purely computational and relates to the issue of coming up with good proposals  $f_{\mathcal{E}, \mathbf{B}}^m$  such that the pseudo-data are sufficiently close to  $\mathbf{Y}^{obs}$ . Due to the vastness of the tree ensemble space, it would be naive to think that one can obtain solid guesses of  $f_0$  purely by sampling from non-informative priors. This is why we call this ABC implementation naive. These considerations lead us to a new data-splitting ABC modification that uses a random portion of the data to train the prior and to generate pseudo-data with more affinity to the left-out observations.

### 3.2. ABC Bayesian Forests

By sampling directly from noninformative priors over tree ensembles  $\pi(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathcal{S})$ , the acceptance rate of the naive ABC can be prohibitively small where huge reference tables would be required to obtain only a few approximate samples from the posterior.

To address this problem, we suggest a sample-splitting approach to come up with draws that are less likely to be rejected by the ABC method. At each ABC iteration, we first draw a random subsample  $\mathcal{I} \subset \{1, \dots, n\}$  of size  $|\mathcal{I}| = s$  with no replacement. Then we split the observed data  $\mathbf{Y}^{(n)}$  into two groups, denoted with  $\mathbf{Y}_{\mathcal{I}}^{(n)}$  and  $\mathbf{Y}_{\mathcal{I}^c}^{(n)}$ , and instead of (6) we consider the marginal likelihood conditionally on  $\mathbf{Y}_{\mathcal{I}}^{(n)}$

$$\pi(\mathbf{Y}^{(n)} | \mathbf{Y}_{\mathcal{I}}^{(n)}, \mathcal{S}) = \int_{(f_{\mathcal{E}, \mathbf{B}}, \sigma^2)} \pi(\mathbf{Y}_{\mathcal{I}^c}^{(n)} | f_{\mathcal{E}, \mathbf{B}}, \sigma^2) d\pi_{\mathcal{I}}(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathcal{S}) \quad (7)$$

where

$$\pi_{\mathcal{I}}(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathcal{S}) = \pi(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathbf{Y}_{\mathcal{I}}^{(n)}, \mathcal{S}). \quad (8)$$

This simple decomposition unfolds new directions for ABC sampling based on data splitting. Instead of using all observations  $\mathbf{Y}^{obs}$  to accept/reject each draw, we set aside a random subset of data  $\mathbf{Y}_{\mathcal{I}^c}^{obs}$  for ABC rejection and use  $\mathbf{Y}_{\mathcal{I}}^{obs}$  to “train the prior”. The key observation is that the samples from the prior  $\pi_{\mathcal{I}}(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathcal{S})$ , i.e. the *posterior*  $\pi(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathbf{Y}_{\mathcal{I}}^{(n)}, \mathcal{S})$ , will have seen a part of the data and will produce more realistic guesses of  $f_0$ . Such guesses are more likely to yield pseudo-data that match  $\mathbf{Y}_{\mathcal{I}^c}^{obs}$  more closely, thereby increasing the acceptance rate of ABC sampling. Note that the acceptance step is based solely on the left-out sample  $\mathbf{Y}_{\mathcal{I}^c}^{obs}$ , not the entire data. Similarly as the naive ABC outlined in the previous section, we first sample the subset  $\mathcal{S}$  from the prior  $\pi(\mathcal{S})$  and then obtain draws from the conditional marginal likelihood under an updated prior  $\pi_{\mathcal{I}}(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathcal{S})$ . This corresponds to an ABC strategy for sampling from  $\pi(\mathcal{S} | \mathbf{Y}_{\mathcal{I}^c}^{(n)})$  under the priors (2) and (8). As will be seen later, this posterior is effective



for assessing variable importance. Moreover, if  $\pi(\mathcal{S})$  is a good proxy for  $\pi(\mathcal{S} | \mathbf{Y}_{\mathcal{I}}^{(n)})$  (when the training set is small relative to the ABC rejection set), this ABC will produce approximate samples from the original target  $\pi(\mathcal{S} | \mathbf{Y}^{(n)})$ .

The idea of using a portion of the data for training the prior and the rest for model selection goes back to at least Good (1950). The most common prescription for choosing training samples in Bayesian analysis is to convert improper priors into proper ones for meaningful model selection with Bayes factors (Lempers, 1971; O’Hagan, 1995). Berger and Pericchi (1996) advocated choosing the training set as small as possible subject to yielding proper posteriors (so called minimal training samples). Berger and Pericchi (2004) argue that data can vary widely in terms of their information content and the use of single minimal training samples can be inadequate/ suboptimal. Since there are many possible training samples, it is natural to average the resulting Bayes factors over the training samples in some fashion. While intrinsic Bayes factors (Berger and Pericchi, 1996) average Bayes factors over all possible minimal training samples, expected posterior priors (Pérez and Berger, 2002) average the prior first. In particular, the empirical expected-posterior prior for model  $\mathcal{S}$  (Ghosh and Samanta, 2002; Pérez and Berger, 2002) writes as

$$\pi(f_{\mathcal{E}, \mathcal{B}}, \sigma^2 | \mathcal{S}) = \frac{1}{L} \sum_{l=1}^L \pi_{\mathcal{I}_l}(f_{\mathcal{E}, \mathcal{B}}, \sigma^2 | \mathcal{S}), \quad (9)$$

where  $\pi_{\mathcal{I}_l}(f_{\mathcal{E}, \mathcal{B}}, \sigma^2 | \mathcal{S})$  was defined in (8) and where  $L$  is the number of all minimal training samples  $\mathcal{I}_l$ . The marginal likelihood under this prior can be then written as (equation (3.5) in Pérez and Berger (2002))  $m(\mathbf{Y}^{(n)} | \mathcal{S}) = \frac{1}{L} \sum_{l=1}^L \pi(\mathbf{Y}^{(n)} | \mathbf{Y}_{\mathcal{I}}^{(n)}, \mathcal{S})$ , where  $\pi(\mathbf{Y}^{(n)} | \mathbf{Y}_{\mathcal{I}}^{(n)}, \mathcal{S})$  was defined in (7). Our ABC analysis with internal data splitting can be thus regarded as arising from the empirical expected posterior prior (9). While the motivation for using training samples in Bayesian analysis has been largely to make improper priors proper, here we use this idea in a different context to increase ABC acceptance rate.

The ABC Bayesian Forests algorithm is formally summarized in Table 1. It starts by splitting the dataset into two subsets at each ( $m^{\text{th}}$ ) iteration:  $\mathbf{Y}_{\mathcal{I}_m}^{\text{obs}}$  for fitting and  $\mathbf{Y}_{\mathcal{I}_m^c}^{\text{obs}}$  for ABC rejection. The algorithm then proceeds by sampling an active set  $\mathcal{S}$  from  $\pi(\mathcal{S})$ . Using the spike-and-slab construction, one can draw Bernoulli indicators  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)'$  where  $\mathbb{P}(\gamma_j = 1 | \theta) = \theta$  for some prior inclusion probability  $\theta \in (0, 1)$  and set  $\mathcal{S}_m = \{j : \gamma_j = 1\}$ . When sparsity is anticipated, one can choose  $\theta$  to be small or to arise from a beta prior  $\mathcal{B}(a, b)$  for some  $a > 0$  and  $b > 0$  (yielding the beta-binomial prior). We discuss other suitable prior model choices in Section 4.

In the (c) step of ABC Bayesian Forests, one obtains a sample from the posterior of  $(f_{\mathcal{E}, \mathcal{B}}, \sigma^2)$ , given  $\mathbf{Y}_{\mathcal{I}_m}^{\text{obs}}$ . For this step, one can leverage existing implementations of Bayesian CART and BART (e.g. the BART R package of McCulloch et al. (2018)). A single draw from the posterior is obtained after a sufficient burn-in. In this vein, one can view ABC Bayesian Forests as a computational envelope around BART to restrict the pool of available variables. The (d) step then consists of predicting the outcome  $\mathbf{Y}_{\mathcal{I}_m^c}^*$  for left-out observations  $\mathbf{x}_i$  using (1) for each  $i \in \mathcal{I}_m^c$ . The last step is ABC rejection based on the discrepancy between  $\mathbf{Y}_{\mathcal{I}_m^c}^*$  and  $\mathbf{Y}_{\mathcal{I}_m^c}^{\text{obs}}$ .

**Algorithm 1 : ABC Bayesian Forests****Data:** Data  $(Y_i^{obs}, \mathbf{x}_i)$  for  $1 \leq i \leq n$ **Result:**  $\pi_j(\epsilon)$  for  $1 \leq j \leq p$  where  $\pi_j(\epsilon) = \widehat{\mathbb{P}}(j \in \mathcal{S}_0 | \mathbf{Y}^{(n)})$ **Set**  $M$ : the number of ABC simulations;  $s$ : the subsample size;  $\epsilon$ : the tolerance threshold;  $m = 0$  the counter**while**  $m \leq M$  **do**(a) **Split** data  $\mathbf{Y}^{obs}$  into  $\mathbf{Y}_{\mathcal{I}_m}^{obs}$  and  $\mathbf{Y}_{\mathcal{I}_m^c}^{obs}$ , where  $\mathcal{I}_m \subset \{1, \dots, n\}$  of size  $|\mathcal{I}_m| = s$  is obtained by sampling with no replacement.(b) **Pick** a subset  $\mathcal{S}_m$  from  $\pi(\mathcal{S})$ .(c) **Sample**  $(f_{\mathcal{E}, \mathbf{B}}^m, \sigma_m^2)$  from  $\pi_{\mathcal{I}_m}(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathcal{S}_m) = \pi(f_{\mathcal{E}, \mathbf{B}}, \sigma^2 | \mathbf{Y}_{\mathcal{I}_m}^{obs}, \mathcal{S}_m)$ .(d) **Generate** pseudo-data  $\mathbf{Y}_{\mathcal{I}_m^c}^*$  by sampling white noise  $\varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_m^2)$  and setting  $Y_i^* = f_{\mathcal{E}, \mathbf{B}}^m(\mathbf{x}_i) + \varepsilon_i$  for each  $i \notin \mathcal{I}_m$ .(e) **Compute** discrepancy  $\epsilon_m = \|\mathbf{Y}_{\mathcal{I}_m^c}^* - \mathbf{Y}_{\mathcal{I}_m^c}^{obs}\|_2$ .**if**  $\epsilon_m < \epsilon$  **then**| Accept  $(\mathcal{S}_m, f_{\mathcal{E}, \mathbf{B}}^m)$  and set  $m = m + 1$ **else**| Reject  $(\mathcal{S}_m, f_{\mathcal{E}, \mathbf{B}}^m)$  and set  $m = m + 1$ **end****end****Compute**  $\pi_j(\epsilon)$  as the proportion of times  $j^{th}$  variable is used in the accepted  $f_{\mathcal{E}, \mathbf{B}}^m$ 's.

For the computation of marginal inclusion probabilities  $\pi_j(\epsilon)$ , one could conceivably report the proportion of ABC accepted samples  $\mathcal{S}_m$  that contain the  $j^{th}$  variable. However,  $\mathcal{S}_m$  is a pool of *available* predictors and not all of them are necessarily used in  $f_{\mathcal{E}, \mathbf{B}}^m$ . Thereby, we report the proportion of ABC accepted samples  $f_{\mathcal{E}, \mathbf{B}}^m$  that use the  $j^{th}$  variable at least once, i.e.

$$\pi_j(\epsilon) = \frac{1}{M(\epsilon)} \sum_{m: \epsilon_m < \epsilon} \mathbb{I}(j \text{ used in } f_{\mathcal{E}, \mathbf{B}}^m), \quad (10)$$

where  $M(\epsilon)$  is the number of accepted ABC samples at  $\epsilon$ . Each tree ensemble  $f_{\mathcal{E}, \mathbf{B}}^m$  thus performs its own variable selection by picking variables from  $\mathcal{S}_m$  rather than from  $\{1, \dots, p\}$ . Limiting the pool of predictors prevents from too many false positives. In addition, the inclusion probabilities (10) do use the training data  $\mathbf{Y}_{\mathcal{I}}^{(n)}$  to shrink and update the subset  $\mathcal{S}$  by leaving out covariates not picked by  $f_{\mathcal{E}, \mathbf{B}}^m$ . In this way, the mechanism for selecting the subsets  $\mathcal{S}$  is not strictly sampling from the prior  $\pi(\mathcal{S})$  but it seizes the information in the training set  $\mathcal{I}$ . In this way,  $\mathcal{S}_m$ 's can be regarded as approximate samples from  $\pi(\mathcal{S} | \mathbf{Y}^{obs})$ . When  $\mathcal{I} = \emptyset$ , we recover the naive ABC as a special case.

**3.2.1. Dynamic ABC**

The estimates of marginal inclusion probabilities  $\pi_j(\epsilon)$  obtained with ABC Bayesian Forests unavoidably depend on the level of approximation accuracy  $\epsilon$ . The acceptance threshold  $\epsilon$  can be difficult to determine in practice, because it has to accommodate random variation of data around  $f_0$  as well as the error when approximating smooth

surfaces  $f_0$  with trees. As  $\epsilon \rightarrow 0$ , the approximations  $\pi_j(\epsilon)$  will be more accurate, but the acceptance rate will be smaller. It is customary to pick  $\epsilon$  as an empirical quantile of  $\epsilon_m$  (Grelaud et al., 2009), keeping only the top few closest samples. Rather than choosing one value  $\epsilon$ , we suggest a dynamic strategy by considering a sequence of decreasing values  $\epsilon_N > \epsilon_{N-1} > \dots > \epsilon_1 > 0$ . By filtering out the ABC samples with stricter thresholds, we track the evolution of each  $\pi_j(\epsilon)$  as  $\epsilon$  gets smaller and smaller. This gives us a dynamic plot that is similar in spirit to the Spike-and-Slab LASSO (Ročková and George, 2018) or EMVS (Ročková and George, 2014) coefficient evolution plots. However, our plots depict approximations to posterior inclusion probabilities rather than coefficient magnitudes. Other strategies for selecting the threshold  $\epsilon$  are discussed in (Sunnaaker et al., 2013; Marin et al., 2012; Csillery et al., 2010).

### 3.3. ABC Bayesian Forests in Action

We demonstrate the usefulness of ABC Bayesian Forests on the benchmark Friedman dataset (Friedman, 1991), where the observations are generated from (1) with  $\sigma = 1$  and

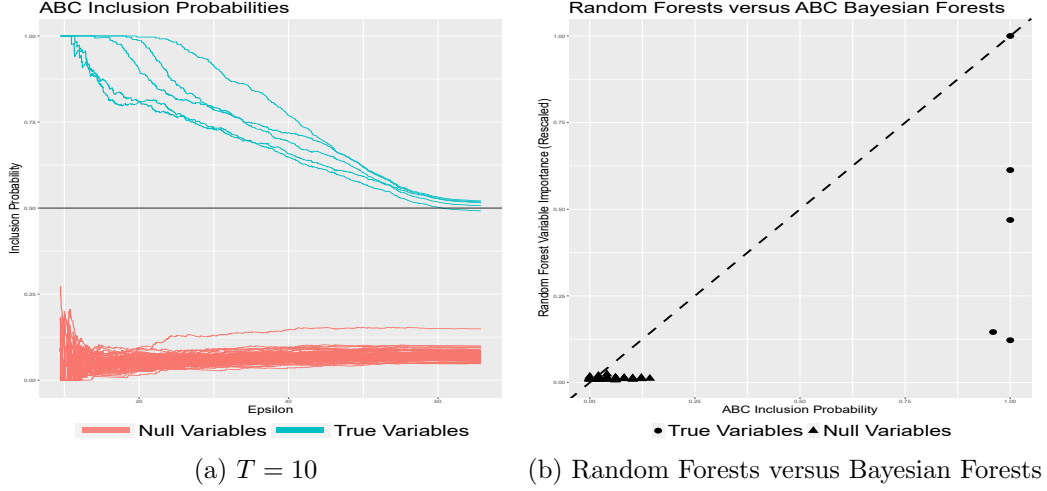
$$f_0(\mathbf{x}_i) = 10 \sin(\pi x_{i1} x_{i2}) + 20 (x_{i3} - 0.5)^2 + 10 x_{i4} + 5 x_{i5}, \quad (11)$$

where  $x_i \in [0, 1]^p$  are *iid* from a uniform distribution on a unit cube. Because the outcome depends on  $x_1, \dots, x_p$ , the predictors  $x_6, \dots, x_p$  are irrelevant, making it more challenging to find  $f_0(\mathbf{x})$ . We begin by illustrating the basic features of ABC Bayesian Forests with  $p = 100$  and  $n = 500$ , assuming the beta-binomial prior  $\pi(\mathcal{S} | \theta)$  with  $\theta \sim \mathcal{B}(1, 1)$  (see Section 3.2). At the  $m^{\text{th}}$  ABC iteration, we draw one posterior sample  $f_{\mathcal{E}, \mathcal{B}}^m$  after 100 burnin iterations using the BART MCMC algorithm (Chipman et al., 2001) with  $T = 10$  trees. We generate  $M = 1000$  ABC samples (with  $s = n/2$ ) and we keep track of variables used in  $f_{\mathcal{E}, \mathcal{B}}^m$ 's to estimate the marginal posterior inclusion probabilities  $\pi_j(\epsilon)$ . It is worth pointing out that unlike MCMC, ABC Bayesian Forests are embarrassingly parallel, making distributed implementations readily available.

Following the dynamic ABC strategy, we plot the estimates of posterior inclusion indicators  $\pi_j(\epsilon)$  as a function of  $\epsilon$  (Figure 1). The true signals are depicted in blue, while the noise covariates are in red. The estimated inclusion probabilities clearly segregate the active and non-active variables, even for large  $\epsilon$  values. This is because BART itself performs variable selection to some degree, where not all variables in  $\mathcal{S}_m$  end up contributing to  $f_{\mathcal{E}, \mathcal{B}}^m$ . For small enough  $\epsilon$ , the inclusion probabilities of true signals eventually cross the 0.5 threshold. Based on the median probability model rule (Barbieri and Berger, 2004), one thereby selects the true model when  $\epsilon$  is sufficiently small. Because the inclusion probabilities get a bit unstable as  $\epsilon$  gets smaller (they are obtained from smaller reference tables), we excluded the 10 smallest  $\epsilon$  values from the plot.

We repeated the experiment with more trees ( $T = 50$ ) and a single tree ( $T = 1$ ). Using more trees, one still gets the separation between signal and noise. However, many more noisy covariates would be included by the MPM rule. This is in accordance with Chipman et al. (2001) who state that BART can over-select with many trees. With a single tree, on the other hand, one may miss some of the low-signal predictors, where deeper trees and more ABC iterations would be needed to obtain a clearer separation.

In this simulation, we observe a curious empirical connection between  $\pi_j(\epsilon)$ , obtained with ABC Bayesian Forests (taking top 5% ABC samples), and rescaled variable im-



**Fig. 1.** (Left) Dynamic ABC plots for evolving inclusion probabilities as  $\epsilon$  gets smaller. (Right) Plot of  $\pi_j(\epsilon)$  obtained with ABC Bayesian Forests ( $\epsilon$  is the 5% quantile of  $\epsilon_m$ 's) and the variable importance measure from Random Forests (rescaled to have a maximum at 1).

portances obtained with Random Forests (RF). From Figure 1(b), we see that the two measures largely agree, separating the signal coefficients (triangles) from the noise coefficients (dots). However, the RF measure is a bit more conservative, yielding smaller normalized importance scores for true signals. While variable importance for RF is yet not understood theoretically, in the next section we provide conditions under which the posterior distribution is consistent for variable selection.

#### 4. Model-Free Variable Selection Consistency

In this section, we develop large sample model selection theory for spike-and-forest priors. As a jumping-off point, we first assume that  $\alpha$  (the regularity of  $f_0$ ) is known, where model selection essentially boils down to finding the active set  $\mathcal{S}_0$ . Later in this section, we investigate *joint* model selection consistency, acknowledging uncertainty about  $\mathcal{S}_0$  and, at the same time, the regularity  $\alpha$ .

Several consistency results for non-parametric regression already exist (Zhu et al., 2015; Yang and Pati, 2017). Comminges and Dalalyan (2012) characterized tight conditions on  $(n, p, q_0)$ , under which it is possible to consistently estimate the sparsity pattern in two regimes. For fixed  $q_0$ , consistency is attainable when  $(\log p)/n \leq c$  for some  $c > 0$ . When  $q_0$  tends to infinity as  $n \rightarrow \infty$ , consistency is achievable when  $c_1 q_0 + \log \log(p/q_0) - \log n \leq c_2$  for some  $c_1, c_2 > 0$ . Throughout this section, we will treat  $q_0$  as fixed and show variable selection consistency when  $q_0 \log p \leq n^{q_0/(2\alpha+q_0)}$ . As an overture to our main result, we start with a simpler case when  $T = 1$  (a single tree) and when  $\alpha$  is known. The full-fledged result for Bayesian forests and unknown  $\alpha$  is presented in Section 4.3. Throughout this section, we will assume  $\sigma^2 = 1$ .

#### 4.1. The Case of Known $\alpha$

Spike-and-forest mixture priors are constructed in two steps by (1) first specifying a conditional prior  $\Pi_{\mathcal{S}}(f)$  on tree (ensemble) functions expressing a qualitative guess on  $f_0$ , and then (2) attaching a prior weight  $\pi(\mathcal{S})$  to each “model” (i.e. subset)  $\mathcal{S}$ . The posterior distribution  $\Pi(f | \mathbf{Y}^{(n)})$  can be viewed as a mixture of individual posteriors for various models  $\mathcal{S}$  with weights given by posterior model probabilities  $\Pi(\mathcal{S} | \mathbf{Y}^{(n)})$ , i.e.

$$\Pi(f | \mathbf{Y}^{(n)}) = \sum_{\mathcal{S}} \Pi(\mathcal{S} | \mathbf{Y}^{(n)}) \Pi_{\mathcal{S}}(f | \mathbf{Y}^{(n)}).$$

Our aim is to establish “model-free” variable selection consistency in the sense that

$$\Pi(\mathcal{S} = \mathcal{S}_0 | \mathbf{Y}^{(n)}) \rightarrow 1 \quad \text{in } \mathbb{P}_{f_0}^{(n)}\text{-probability as } n \rightarrow \infty,$$

where  $\mathbb{P}_{f_0}^{(n)}$  is the distribution of  $\mathbf{Y}^{(n)}$  under (1). The adjective “model-free” merely refers to the fact that we are selecting subsets in a non-parametric regression environment without necessarily committing to a linear model. We start by defining the model index set  $\mathbf{\Gamma} = \{\mathcal{S} : \mathcal{S} \subseteq \{1, \dots, p\}\}$ , consisting of all  $2^p$  variable subsets, and we partition it into (a) the true model  $\mathcal{S}_0$ , (b) models that *overfit*  $\mathbf{\Gamma}_{\mathcal{S} \supset \mathcal{S}_0}$  (i.e. supersets of the true subset  $\mathcal{S}_0$ ) and (c) models that *underfit*  $\mathbf{\Gamma}_{\mathcal{S} \not\supset \mathcal{S}_0}$  (i.e. models that miss at least one active covariate). Each model  $\mathcal{S} \in \mathbf{\Gamma}$  is accompanied by a convergence rate  $\varepsilon_{n,\mathcal{S}}$  that reflects the inherent difficulty of the estimation problem. For each model  $\mathcal{S}$  of size  $|\mathcal{S}|$ , we define

$$\varepsilon_{n,\mathcal{S}} = C_{\varepsilon} n^{-\alpha/(2\alpha+|\mathcal{S}|)} \sqrt{\log n} \quad \text{for some } C_{\varepsilon} > 0, \quad (12)$$

the  $\|\cdot\|_n$ -near-minimax rate of estimation of a  $|\mathcal{S}|$ -dimensional  $\alpha$ -smooth function.

##### 4.1.1. Prior Specification

Prior distribution on the model index  $\Pi(\mathcal{S})$  has to be chosen carefully for model selection consistency to hold when  $p > n$  (Moreno et al., 2015). Traditional spike-and-slab priors introduce  $\Pi(\mathcal{S})$  through a prior inclusion probability  $\theta = \Pi(i \in \mathcal{S}_0 | \theta)$ , independently for each  $i = 1, \dots, p$ . This prior mixing weight is often endowed with a prior, such as the uniform prior  $\pi(\theta) = \mathcal{B}(1, 1)$  (Scott and Berger, 2010), yielding a uniform prior on the model size, or the “complexity prior”  $\pi(\theta) = \mathcal{B}(1, p^c)$  for  $c > 2$  (Castillo and van der Vaart, 2012), yielding an exponentially decaying prior on the model size. We propose a different approach, directly assigning a prior on model weights through

$$\pi(\mathcal{S}) \propto e^{-C(n^{|\mathcal{S}|/(2\alpha+|\mathcal{S}|)} \log n \vee |\mathcal{S}| \log p)} \quad (13)$$

where  $C > 0$  is a suitably large constant. When  $|\mathcal{S}| \log p \leq n^{|\mathcal{S}|/(2\alpha+|\mathcal{S}|)}$ , this prior is proportional to  $e^{-C/C_{\varepsilon}^2 n \varepsilon_{n,\mathcal{S}}^2}$  and, as such, it puts more mass on models that yield faster rates convergence (similarly as in Lember and van der Vaart (2007)). When  $|\mathcal{S}| \log p > n^{|\mathcal{S}|/(2\alpha+|\mathcal{S}|)} \log n$ , the implied prior on the effective dimensionality  $\pi(|\mathcal{S}|) = \binom{p}{|\mathcal{S}|} \pi(\mathcal{S})$  will be exponentially decaying in the sense that  $\pi(|\mathcal{S}|) \lesssim e^{-(C-1)|\mathcal{S}| \log p}$  for  $C > 1$ . It was recently noted by Castillo and Mismar (2018) that the complexity prior “penalizes slightly more than necessary”. With our prior specification (13), however, the exponential decay kicks in *only* when  $|\mathcal{S}|$  is sufficiently large.

Assuming that the level of smoothness  $\alpha$  is known, the optimal number of steps (i.e. tree bottom leaves  $K$ ) needed to achieve the rate-optimal performance for estimating  $f_0$  should be of the order  $n^{q_0/(2\alpha+q_0)} = 1/C_\varepsilon^2 n \varepsilon_{n, \mathcal{S}_0}^2 / \log n$  (Ročková and van der Pas, 2017). For our toy setup with a known  $\alpha$ , we thus assume a point-mass prior on  $K$  with an atom near the optimal number of steps for each given  $\mathcal{S}$ , i.e.

$$\pi(K | \mathcal{S}) = \mathbb{I}[K = K_{\mathcal{S}}], \quad \text{where} \quad K_{\mathcal{S}} = \lfloor C_K / C_\varepsilon^2 n \varepsilon_{n, \mathcal{S}}^2 / \log n \rfloor \quad (14)$$

for some  $C_K > 0$  such that  $K_{\mathcal{S}_0} = 2^{q_0 s}$  for some  $s \in \mathbb{N}$ . In Section 4.2, we allow for more flexible trees with variable sizes.

#### 4.1.2. Identifiability

The active variables ought to be sufficiently relevant in order to make their identification possible. To this end, we introduce a non-parametric signal strength assumption, making sure that  $f_0$  is not too flat in active directions (Yang and Pati, 2017; Comminges and Dalalyan, 2012).

We first introduce the notion of an approximation gap. For any given model  $\mathcal{S}$ , we denote with  $\mathcal{F}_{\mathcal{S}}$  a set of approximating functions (only single trees  $f_{\mathcal{T}, \beta}$  with  $K_{\mathcal{S}}$  leaves for now) and define the approximation gap as follows:

$$\delta_n^{\mathcal{S}} \equiv \inf_{f_{\mathcal{T}, \beta} \in \mathcal{F}_{\mathcal{S}}} \|f_0 - f_{\mathcal{T}, \beta}\|_n = \|f_0 - f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}}\|_n, \quad (15)$$

where  $f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}}$  is the  $\|\cdot\|_n$ -projection of  $f_0$  onto  $\mathcal{F}_{\mathcal{S}}$ . For identifiability of  $\mathcal{S}_0$ , we require that those models that miss one of the active covariates have a large separation gap.

DEFINITION 4.1. (*Identifiability*) We say that  $\mathcal{S}_0$  is  $(f_0, \varepsilon)$ -identifiable if, for some  $M > 0$ ,

$$\inf_{i \in \mathcal{S}_0} \delta_n^{\mathcal{S}_0 \setminus i} > 2M\varepsilon. \quad (16)$$

We provide a more intuitive explanation of (16) in terms of directional variability of  $f_0$ . The best approximating tree  $f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}}$  can be written as

$$f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}}(\mathbf{x}) = \sum_{k=1}^{K_{\mathcal{S}}} \mathbb{I}(\mathbf{x} \in \widehat{\Omega}_k^{\mathcal{S}}) \widehat{\beta}_k \quad \text{with} \quad \widehat{\beta}_k = \bar{f}_0(\widehat{\Omega}_k^{\mathcal{S}}) \equiv \frac{1}{n(\widehat{\Omega}_k^{\mathcal{S}})} \sum_{\mathbf{x}_i \in \widehat{\Omega}_k^{\mathcal{S}}} f_0(\mathbf{x}_i),$$

where  $\widehat{\mathcal{T}} = \{\widehat{\Omega}_k^{\mathcal{S}}\}_{k=1}^{K_{\mathcal{S}}}$  is the tree-shaped partition of the  $\|\cdot\|_n$ -projection of  $f_0$  defined in (15) with  $K_{\mathcal{S}}$  leaves and where  $n(\widehat{\Omega}_k^{\mathcal{S}}) = \sum_{i=1}^n \mathbb{I}(\mathbf{x}_i \in \widehat{\Omega}_k^{\mathcal{S}}) \equiv n \mu(\widehat{\Omega}_k^{\mathcal{S}})$ . The separation gap in (15) can be then re-written as

$$\delta_n^{\mathcal{S}} = \sqrt{\sum_{k=1}^{K_{\mathcal{S}}} \mu(\widehat{\Omega}_k^{\mathcal{S}}) V[f_0 | \widehat{\Omega}_k^{\mathcal{S}}]},$$

where

$$V[f_0 | \widehat{\Omega}_k^{\mathcal{S}}] \equiv \frac{1}{n(\widehat{\Omega}_k^{\mathcal{S}})} \sum_{\mathbf{x}_i \in \widehat{\Omega}_k^{\mathcal{S}}} (f_0(\mathbf{x}_i) - \bar{f}_0(\widehat{\Omega}_k^{\mathcal{S}}))^2$$



is the local variability of  $f_0$  inside  $\widehat{\Omega}_k^{\mathcal{S}}$ . Given this characterization, (16) will be satisfied, for instance, when variability of  $f_0$  inside best approximating cells that miss an active direction is too large, i.e.  $\inf_{i \in \mathcal{S}_0} \inf_k V[f_0 | \widehat{\Omega}_k^{\mathcal{S}_0 \setminus i}] > 4M^2 \varepsilon^2$ .

Our identifiability condition is a theoretical assumption on  $f_0$  which indicates how large signal in each direction should be in order to be capturable. It generalizes the more traditional sufficient “beta-min conditions” (Castillo et al., 2015; Zhao and Yu, 2006) for variable selection consistency (see Remark 4.1). Here, we gauge the amount of signal in terms of local variation in cells that *do not split* on an active covariate. Intuitively, if we do not split on  $i \in \mathcal{S}_0$ , the “variation” of  $f_0$  inside the cells of the best tree we can get without  $i$  will be too large. The following example links our identifiability assumption with beta-min conditions.

EXAMPLE 4.1. Assume for now that  $p = 2$  and that  $f_0$  is linear, i.e.

$$f_0(\mathbf{x}_i) = a + bx_{i1} + cx_{i2}.$$

Moreover, assume that  $n = 16$  predictor observations are located on a regular grid  $\mathcal{X} = \{k/4 : 1 \leq k \leq 4\} \times \{j/4 : 1 \leq j \leq 4\}$ , where  $\times$  denotes the Cartesian product. Suppose  $\mathcal{S}_0 = \{1, 2\}$  and set  $\mathcal{S} = \mathcal{S}_0 \setminus \{2\} = \{1\}$  and  $K_{\mathcal{S}} = 2$ . It can be verified that the partition  $\widehat{\mathcal{T}}$  of the best approximating tree that does not split on the covariate  $x_2$  consists of two rectangles  $\widehat{\Omega}_1^{\mathcal{S}} = [0, 1/2) \times [0, 1]$  and  $\widehat{\Omega}_2^{\mathcal{S}} = [1/2, 1] \times [0, 1]$ . Then we have

$$\bar{f}_0(\widehat{\Omega}_1^{\mathcal{S}}) = a + \frac{3}{2} \left( \frac{b}{4} \right) + \frac{5}{2} \left( \frac{c}{4} \right) \quad \text{and} \quad \bar{f}_0(\widehat{\Omega}_2^{\mathcal{S}}) = a + \frac{7}{2} \left( \frac{b}{4} \right) + \frac{5}{2} \left( \frac{c}{4} \right)$$

and thereby

$$(\delta_m^{\mathcal{S}})^2 = V(f_0 | \widehat{\Omega}_1^{\mathcal{S}}) = V(f_0 | \widehat{\Omega}_2^{\mathcal{S}}) = \frac{1}{4} \frac{b^2}{16} + \frac{5}{4} \frac{c^2}{16}. \quad (17)$$

From the expression (17) we can immediately see the connection to the beta-min conditions. When the signal in the direction of  $x_2$  is large enough, i.e.  $c > 16/\sqrt{5}M\varepsilon$ , our identifiability condition will be satisfied.

The second sufficient condition needed for methods such as the LASSO to fully recover  $\mathcal{S}_0$  is “irrepresentability” (Zhao and Yu, 2006; Van De Geer and Bühlmann, 2009). This condition restricts the amount of correlation between (active and non-active) covariates by imposing a regularization constraint on the magnitudes of regression coefficients of the inactive predictors onto the active ones. Here, we generalize the notion of irrepresentability to the non-parametric setup. Consider an underfitting model  $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2 \not\supseteq \mathcal{S}_0$ , where  $\mathcal{S}_1 \subset \mathcal{S}_0$  are true positives and  $\mathcal{S}_2$  is a possibly empty set of false positives, i.e.  $\mathcal{S}_2 \cap \mathcal{S}_0 = \emptyset$ . We define

$$\rho_n^{\mathcal{S}} \equiv \frac{1}{n} \sum_{i=1}^n [f_0(\mathbf{x}_i) - f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}_1}(\mathbf{x}_i)] [f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}}(\mathbf{x}_i) - f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}_1}(\mathbf{x}_i)], \quad (18)$$

the sample covariance between the surplus signals in  $f_0$  and  $f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}}$  obtained by removing the effect of  $f_{\widehat{\mathcal{T}}, \widehat{\beta}}^{\mathcal{S}_1}$ . This quantity will be large if noise covariates inside  $\mathcal{S}_2$  can compensate for the missed true covariates in  $\mathcal{S}_0 \setminus \mathcal{S}_1$ , i.e. when the true and fake covariates are strongly correlated. To obviate this substitution effect, we introduce the following nonparametric

“irrepresentability” condition. Similarly as in [Zhao and Yu \(2006\)](#), we require that “the total amount of an irrelevant covariate represented by the covariates in the true model” is small.

**DEFINITION 4.2.** (*Irrepresentability*) *We say that  $\varepsilon$ -irrepresentability holds for  $f_0$  and  $\mathcal{S}_0$  if, for some  $M > 0$ , we have  $\sup_{\mathcal{S} \not\supset \mathcal{S}_0} |\rho_n^{\mathcal{S}}| < \frac{M}{2}\varepsilon$ , where  $\rho_n^{\mathcal{S}}$  was defined in (18).*

It follows from Lemma [S.1.2](#) (Appendix) that under the irrepresentability and identifiability conditions (Definition [4.1](#) and [4.2](#)), we obtain

$$\inf_{\mathcal{S} \not\supset \mathcal{S}_0} \inf_{f_{\mathcal{T}, \beta} \in \mathcal{F}_{\mathcal{S}}} \|f_{\mathcal{T}, \beta} - f_0\|_n > M\varepsilon. \quad (19)$$

This condition essentially states that *all* models that miss *at least one* active covariate (i.e. not only subsets of the true model) have a large separation gap.

The following theorem characterizes variable selection consistency of spike-and-tree posterior distributions. Namely, the posterior distribution over the model index is shown to concentrate on the true model  $\mathcal{S}_0$ . One additional assumption is needed to make sure that the (fixed) design  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is sufficiently regular. [Ročková and van der Pas \(2017\)](#) define the notion of a fixed  $\mathcal{S}_0$ -regular design in terms of cell diameters of a  $k$ - $d$  tree partition (Definition 3.3). This assumption essentially excludes outliers, making sure that the data cloud is spread evenly in active directions (while permitting correlation between covariates).

**THEOREM 4.1.** *Assume  $f_0 \in \mathcal{H}_p^\alpha \cap \mathcal{C}(\mathcal{S}_0)$  for some  $\alpha \in (0, 1]$  and  $\mathcal{S}_0 \subset \{1, \dots, p\}$  with  $q_0 = |\mathcal{S}_0|$  and  $\|f_0\|_\infty \lesssim B$ . Denote with  $\tilde{\varepsilon}_n = C_\varepsilon n^{-\alpha/(2\alpha+q_n)} \sqrt{\log n}$ , where  $q_n = C_q [n \varepsilon_{n, \mathcal{S}_0}^2 / \log p]$  for some  $C_q > 0$ , and assume  $q_0 \log p \leq n^{q_0/(2\alpha+q_0)}$  with  $2 \leq q_0 = \mathcal{O}(1)$  as  $n \rightarrow \infty$ . Assume that (a)  $\mathcal{S}_0$  is  $(f_0, \tilde{\varepsilon}_n)$ -identifiable, (b)  $\tilde{\varepsilon}_n$ -irrepresentability holds and that (c) the design  $\mathcal{X}$  is  $\mathcal{S}_0$ -regular. Under the spike-and-tree prior comprising (with  $T = 1$ ) [\(4\)](#), [\(5\)](#), [\(13\)](#) with  $C > 2$  and [\(14\)](#), we have*

$$\Pi[\mathcal{S} = \mathcal{S}_0 \mid \mathbf{Y}^{(n)}] \rightarrow 1 \quad \text{in } \mathbb{P}_{f_0}^{(n)}\text{-probability as } n \rightarrow \infty.$$

**PROOF.** Section [S.1.1](#)

**REMARK 4.1.** *The assumption of  $(f_0, \tilde{\varepsilon}_n)$ -identifiability pertains to the more traditional sufficient beta-min conditions for variable selection consistency in sparse high-dimensional models. For example, [Castillo et al. \(2015\)](#) in their Corollary 1 require that  $\min_{i \in \mathcal{S}_0} |\beta_i^0| \geq M \sqrt{\frac{q_0 \log p}{n}}$ , for some “large enough constant”  $M > 0$  that depends on the compatibility number (see e.g. Definition 2.1 in [Castillo et al. \(2015\)](#) of the design matrix  $X$  (rescaled to have an  $\|\cdot\|_2$  norm  $\sqrt{n}$ ). Our identifiability threshold also depends on the rate of convergence  $\varepsilon_n$  (similarly as in [Castillo et al. \(2015\)](#)). However, unlike in the linear models we measure the signal strength in a non-parametric way. Lastly, note that the identifiability gap  $\tilde{\varepsilon}_n$  in Theorem [4.1](#) is a bit larger than the near-minimax rate  $\varepsilon_{n, \mathcal{S}_0}$ . This requirement will be relaxed in the next section, where  $\alpha$  will be treated as unknown.*

For *iid* models, Ghosal et al. (2008) considered the problem of nonparametric Bayesian model selection and averaging and characterized conditions under which the posterior achieves adaptive rates of convergence. The authors also study the posterior distribution of the model index, showing that it puts a negligible weight on models that are bigger than the optimal one. Yang and Pati (2017) characterized similar conditions for the non-*iid* case, see Section S.1.1 for more details.

REMARK 4.2. (*Theory for ABC*) It is worth pointing out that Theorem 4.1 is obtained for the actual posterior  $\pi(\mathcal{S} | \mathbf{Y}^{(n)})$ , not the ABC posterior. Theory for ABC recently started emerging with the first results focussing on ABC bias (Barber et al., 2015), consistency and asymptotic normality (Martin et al., 2014; Frazier et al., 2018, 2020) and on convergence of the posterior mean (Li and Fearnhead, 2018). For our non-parametric regression scenario, we can conclude (variable selection) consistency for ABC Bayesian forests under the assumption that the residual variance  $\sigma^2$  decreases with the sample size (as is typical in the Gaussian sequence model). In particular, Theorem S.1.1 in Supplemental Materials (Section S.1.4) shows that the ABC posterior concentrates at the rate  $\lambda_n = 4\epsilon_n^T/3 + 1/\sqrt{n}$ , where  $\epsilon_n^T = \sqrt{2\log n/n}$  is the ABC tolerance level. This result implies that the ABC posterior will not reward underfitting model as long as our identifiability and irrepresentability conditions are satisfied with  $\varepsilon = \lambda_n$ . Regarding over-fitting models, an ABC analogue of Lemma 1.1 (Section 1.1.2 in Supplemental Materials) implies that the ABC posterior probability of over-fitting models goes to zero, which concludes variable selection consistency of a (naive) ABC method. These considerations can be extended to ABC Bayesian Forests with data splitting using the empirical expected posterior prior justification in (9). More details are in Supplemental Materials (Section S.1.4).

REMARK 4.3. (*Consistency of the Median Probability Model*) In Section 3.3, we used the median probability model rule which may not be the same as the highest-posterior model whose consistency we have shown in Theorem 4.1. However, even when  $p \rightarrow \infty$  it can be verified (as in Corollary 4.1 in Narisetty and He (2014)) that the median probability model is also consistent under the same assumptions as Theorem 4.1. In particular,  $\mathbb{P}_{f_0}^{(n)}[\bigcap_{i=1}^p E_i] \rightarrow 1$  as  $n \rightarrow \infty$  where  $E_i = \{\Pi(\gamma_i = \gamma_i^0 | \mathbf{Y}^{(n)}) > 0.5\}$  and where  $\gamma_i = \mathbb{I}(i \in \mathcal{S})$  are binary inclusion indicators and  $\gamma_i^0 = \mathbb{I}(i \in \mathcal{S}_0)$ .

#### 4.2. The Case of Unknown $\alpha$

The fact that the level  $\alpha$  has to be known for the consistency to hold makes the result in Theorem 4.1 somewhat theoretical. In this section, we provide a joint consistency result for the unknown regularity level  $K$  and, at the same time, the unknown subset  $\mathcal{S}_0$ . Finding the optimal regularity level  $K$ , given  $\mathcal{S}_0$ , is a model selection problem of independent interest (Lafferty and Wasserman, 2001). Here, we acknowledge uncertainty about both  $K$  and  $\mathcal{S}_0$  by assigning a joint prior distribution on  $(K, \mathcal{S})$ . Namely, we consider an analogue of (13), where  $n^{|\mathcal{S}|/(2\alpha+|\mathcal{S}|)}$  is now replaced with  $K \log n$  (according to (14)), i.e.

$$\pi(K, \mathcal{S}) \propto e^{-C(K \log n \vee |\mathcal{S}| \log p)} \quad \text{for } 1 \leq K \leq n \quad \text{and } \mathcal{S} \subseteq \{1, \dots, p\}. \quad (20)$$

This prior penalizes models with too many splits or too many covariates. We now regard each model as a *pair of indices*  $(K, \mathcal{S})$ , where the “true” model is characterized by  $\mathbf{\Gamma}_0 = (K_{\mathcal{S}_0}, \mathcal{S}_0)$  with  $K_{\mathcal{S}_0}$  defined in (14). Again, we partition the model index set  $\mathbf{\Gamma} = \{(K, \mathcal{S}) : \mathcal{S} \subseteq \{1, \dots, p\}, 1 \leq K \leq n\}$  into (a) the true model  $\mathbf{\Gamma}_0$ , (b) models that underfit  $\mathbf{\Gamma}_{\{\mathcal{S} \not\supseteq \mathcal{S}_0\} \cup \{K < K_{\mathcal{S}_0}\}}$  (i.e. miss at least one covariate or use less than the optimal number of splits), and (c) models that overfit  $\mathbf{\Gamma}_{\{\mathcal{S} \supset \mathcal{S}_0\} \cap \{K \geq K_{\mathcal{S}_0}\}}$  (i.e. use too many variables and splits).

We combine the identifiability and irrepresentability conditions into one as follows:

$$\inf_{\{\mathcal{S} \not\supseteq \mathcal{S}_0\} \cup \{K < K_{\mathcal{S}_0}\}} \inf_{f_{\mathcal{T}, \beta} \in \mathcal{F}_{\mathcal{S}}(K)} \|f_{\mathcal{T}, \beta} - f_0\|_n > M \varepsilon_{n, \mathcal{S}_0} \quad (21)$$

for some  $M > 1$ , where  $\mathcal{F}_{\mathcal{S}}(K)$  consists of all trees with  $K$  bottom leaves and splitting variables  $\mathcal{S}$ . This condition is an analogue of (19), essentially stating that one cannot approximate  $f_0$  with an error smaller than a multiple of the near-minimax rate using underfitting models.

**THEOREM 4.2.** *Assume  $f_0 \in \mathcal{H}_p^\alpha \cap \mathcal{C}(\mathcal{S}_0)$  for some  $\alpha \in (0, 1]$  and  $\mathcal{S}_0 \subset \{1, \dots, p\}$  such that  $|\mathcal{S}_0| = q_0$  and  $\|f_0\|_\infty \lesssim B$ . Assume  $q_0 \log p \leq n^{q_0/(2\alpha+q_0)}$  and  $2 \leq q_0 = \mathcal{O}(1)$  as  $n \rightarrow \infty$ . Furthermore, assume that the design  $\mathcal{X}$  is  $\mathcal{S}_0$ -regular and that (21) holds. Under the spike-and-tree prior comprising (with  $T = 1$ ) (4), (5) and (20) for  $C > 3$ , we have*

$$\Pi \left[ \{\mathcal{S} = \mathcal{S}_0\} \cap \{K_{\mathcal{S}_0} \leq K \leq K_n\} \mid \mathbf{Y}^{(n)} \right] \rightarrow 1 \quad \text{in } \mathbb{P}_{f_0}^{(n)}\text{-probability as } n \rightarrow \infty,$$

where  $K_{\mathcal{S}_0}$  was defined in (14) and  $K_n = \lceil \bar{C} n \varepsilon_{n, \mathcal{S}_0}^2 / \log n \rceil$  for some  $\bar{C} > C_K / C_\varepsilon^2$ .

**PROOF.** Section S.1.2

Note that both  $K_{\mathcal{S}_0}$  and  $K_n$  are of the same (optimal) order, where the marginal posterior distribution  $\Pi(K \mid \mathbf{Y}^{(n)})$  squeezes inside these two quantities as  $n \rightarrow \infty$ . [Lafferty and Wasserman \(2001\)](#) provide a similar result for their RODEO method, without the variable selection consistency part. [Yang and Pati \(2017\)](#) also provide a similar result for Gaussian processes, without the regularity selection consistency part. Here, we characterize *joint* consistency for both subset and regularity model selection.

### 4.3. Variable Selection Consistency with Bayesian Forests

Finally, we provide a variant of Theorem 4.2 for tree ensembles. Each Bayesian forest (i.e. additive regression tree) model is characterized by a triplet  $(\mathcal{S}, T, \mathbf{K})$ , where  $\mathcal{S}$  is the active variable subset,  $T \in \mathbb{N}$  is the number of trees and  $\mathbf{K} = (K^1, \dots, K^T)' \in \mathbb{N}^T$  is a vector of the bottom leaf counts for the  $T$  trees. Rate-optimality of Bayesian forests can be achieved for a wide variety of priors, ranging from many weak learners (large  $T$  and small  $K^t$ 's) to a few strong learners (small  $T$  and large  $K^t$ 's) ([Ročková and van der Pas, 2017](#)). The optimality requirement is that the *total* number of leaves in the ensemble  $\sum_{t=1}^T K^t$  behaves like  $K_{\mathcal{S}_0}$ , defined earlier in (14).

We thereby define models in terms of equivalence classes rather than individual triplets  $(\mathcal{S}, T, \mathbf{K})$ . We construct each equivalence class  $E(Z)$  by combining ensembles

with the same number  $Z$  of total leaves, i.e.

$$E(Z) = \bigcup_{T=1}^{\min\{Z,n\}} \left\{ \mathbf{K} \in \mathbb{N}^T : \sum_{t=1}^T K^t = Z \right\}. \quad (22)$$

The cardinality of  $E(Z)$ , denoted with  $\Delta(E(Z))$ , satisfies  $\Delta(E(Z)) \leq Z!p(Z)$ , where  $p(Z)$  is the partitioning number (i.e. the number of ways one can write  $Z$  as a sum of positive integers). The “true” model  $\Gamma_0 = (\mathcal{S}_0, E(K_{\mathcal{S}_0}))$  consists of an equivalence class of forests that split on variables inside  $\mathcal{S}_0$  with a total number of  $K_{\mathcal{S}_0}$  leaves. Similarly as before, we define underfitting model classes  $\Gamma_{\{\mathcal{S} \not\supset \mathcal{S}_0\} \cup \{E(Z): Z < K_{\mathcal{S}_0}\}}$  and overfitting model classes  $\Gamma_{\{\mathcal{S} \supset \mathcal{S}_0\} \cap \{E(Z): Z \geq K_{\mathcal{S}_0}\}}$ . Regarding the prior on  $T$ , similarly as Ročková and van der Pas (2017), we consider

$$\pi(T) \propto e^{-C_T T}, \quad T = 1, \dots, n, \quad \text{for } C_T > 0. \quad (23)$$

Given  $T$ , we assign a joint prior over  $\mathcal{S}_0$  and  $\mathbf{K} \in \mathbb{N}^T$  as follows:

$$\pi(\mathcal{S}, \mathbf{K} | T) \propto e^{-C \max\{|\mathcal{S}| \log p; \sum_{t=1}^T K^t \log n\}} \quad \text{for } C > 1. \quad (24)$$

We conclude this section with a model selection consistency result for Bayesian forests under the following identifiability condition

$$\inf_{\{\mathcal{S} \not\supset \mathcal{S}_0\} \cup \{E(Z): Z < K_{\mathcal{S}_0}\}} \inf_{f_{\mathcal{E}, \mathbf{B}} \in \mathcal{F}_{\mathcal{S}}(\mathbf{K})} \|f_{\mathcal{E}, \mathbf{B}} - f_0\|_n > M \varepsilon_{n, \mathcal{S}_0}, \quad (25)$$

where  $\mathcal{F}_{\mathcal{S}}(\mathbf{K})$  denotes all forests  $f_{\mathcal{E}, \mathbf{B}}$  that split on variables  $\mathcal{S}$  and consist of  $T$  trees with  $\mathbf{K} = (K^1, \dots, K^T)'$  bottom leaves.

**THEOREM 4.3.** *Assume  $f_0 \in \mathcal{H}_p^\alpha \cap \mathcal{C}(\mathcal{S}_0)$  for some  $\alpha \in (0, 1]$  and  $\mathcal{S}_0 \subset \{1, \dots, p\}$  such that  $|\mathcal{S}_0| = q_0$  and  $\|f_0\|_\infty \lesssim B$ . Assume  $q_0 \log p \leq n^{q_0/(2\alpha+q_0)}$ , where  $2 \leq q_0 = \mathcal{O}(1)$  as  $n \rightarrow \infty$ . Furthermore, assume that the design is  $\mathcal{S}_0$ -regular and that (25) holds. Under the spike-and-forest prior comprising (4), (5), (23) and (24), we have*

$$\mathbb{P} \left[ \{\mathcal{S} = \mathcal{S}_0\} \cap \left\{ K_{\mathcal{S}_0} \leq \sum_{t=1}^T K^t \leq K_n \right\} \mid \mathbf{Y}^{(n)} \right] \rightarrow 1 \quad \text{in } \mathbb{P}_{f_0}^{(n)}\text{-probability as } n \rightarrow \infty,$$

where  $K_{\mathcal{S}_0}$  was defined in (14) and  $K_n = \lceil \bar{C} n \varepsilon_{n, \mathcal{S}}^2 / \log n \rceil$  for some  $\bar{C} > C_K / C_\varepsilon^2$ .

PROOF. Section S.1.3

## 5. Simulation Study

We evaluate the performance of ABC Bayesian Forests on simulated data. We consider the following performance criteria: Precision =  $1 - \text{FDP} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ , Power =  $\frac{\text{TP}}{\text{TP} + \text{FN}}$  (defined as the proportion of true signals discovered as such), Hamming Distance (HD) =  $\text{FP} + \text{FN}$  (where FP and FN denotes the number of false positives and false negatives, respectively) and the area under the ROC curve (AUC). Traditionally, AUC

assesses how well a classification method can differentiate between two classes in the absence of a clear decision boundary. We use this criterion to assess variable importance since many of the considered selection methods are based on an importance measure and, as such, do not have a clear decision boundary.

The synthetic data are generated from the model (1), where  $\mathbf{x}_i$ 's for  $i = 1, \dots, n$  are drawn independently from  $N_p(0, \Sigma)$  with  $\Sigma = (\rho_{ij})_{i,j=1}^{p,p}$ . We make our comparisons under different combinations of  $f_0$ ,  $\sigma$  and  $\Sigma$ . In particular, we consider a relatively large noise level with  $\sigma = 5$  ( $\sigma = \sqrt{5}$  for the linear setup) and

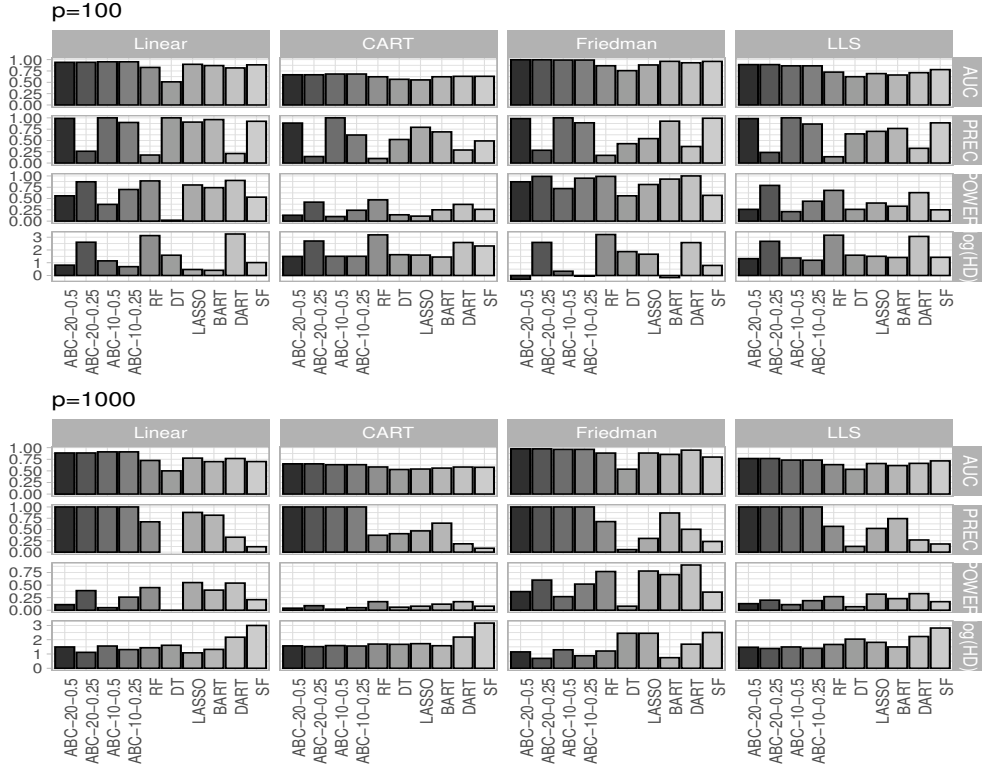
- (a) medium equi-correlation  $\rho_{ij} = 0.5$  for  $i \neq j$  with  $\rho_{ii} = 1$ ,
- (b) high auto-correlation  $\rho_{ij} = 0.9^{|i-j|}$ .

Regarding the mean function  $f_0$ , we consider four choices: (1) a linear setup with  $f_0(\mathbf{x}_i) = x_{i1} + 2x_{i2} + 3x_{i3} - 2x_{i4} - x_{i5}$ ; (2) the Friedman setup as described in (11); (3) a CART (tree-based) function  $f_0(\mathbf{x}_i)$  generated from the first 5 covariates using the `rpart` function in R; (4) a simulated example from Liang et al. (2018) (denoted with LLS hereafter) with  $f_0(\mathbf{x}_i) = \frac{10x_{i2}}{1+x_{i1}^2} + 5 \sin(x_{i3}x_{i4} + 2x_{i5})$ . For the auto-correlation case, we permuted the covariates so that signals are not next to each other.

For each combination of settings, we repeat our simulation over 20 different datasets assuming  $n = 500$  and  $p \in \{100, 1000\}$ . We compare ABC Bayesian Forests with Random Forests (RF), Dynamic Trees (DT) of Taddy et al. (2011b), BART (Chipman et al., 2010), DART of Linero (2018), LASSO and Spike-and-Forests (the MCMC counterpart of ABC Bayesian Forests outlined in Section S.3 of the Supplemental Materials). ABC Bayesian Forests are trained with  $M = 1000$  ABC samples, where only a fraction of ABC samples (top 10%) are kept in the reference table. The prior  $\pi(\mathcal{S})$  is the usual beta-binomial prior with  $\theta \sim \mathcal{B}(1, 1)$ . Inside each ABC step, we sample a subset of size  $s = n/2$  and draw a tree ensemble using the default Bayesian CART prior (Chipman et al., 1998) and  $T \in \{10, 20\}$  trees. For each ABC sample, we draw the last BART sample after  $B = 200$  burnin MCMC iterations. A sensitivity analysis to the choice  $s, T, B$  and  $M$  is reported in the Supplemental Materials (Section 4). Two versions of BART (without ABC) were deployed using the R package BART: (1) the standard BART from Chipman et al. (2010) with  $T = 20$  (as recommended in Bleich et al. (2014)), and (2) the sparse version DART of Linero (2018) with a Dirichlet prior (`sparse=TRUE`, `a=0.5`, `b=1`) with  $T = 200$ . Both versions are run with 10 000 MCMC samples after 10 000 burnin. For LASSO, we use the `glmnet` package in R (Friedman et al., 2010) using the 1-se rule to select the penalty  $\lambda$ . For Random Forests, we deploy the `randomForest` package in R (Liaw and Wiener, 2002) using the default number of 500 trees where variable importance is based on the difference in predictions (with and without each covariate) in out-of-bag samples.

To select variables with random forests, there are at least three commonly used strategies: (1) Recursive Feature Elimination (RFE) implemented in the `caret` package with 5-fold cross-validation (as suggested in Linero (2018)); (2) truncating importance at the  $1 - \alpha$  quantile of a standard normal distribution (as suggested by Breiman and Cutler (2013)); (3) truncating importance at the Bonferroni-corrected  $(1 - \alpha/p)$  quantile of a standard normal distribution (Bleich et al., 2014). We report the third method, which was seen to perform the best. For BART and DART, we select those variables



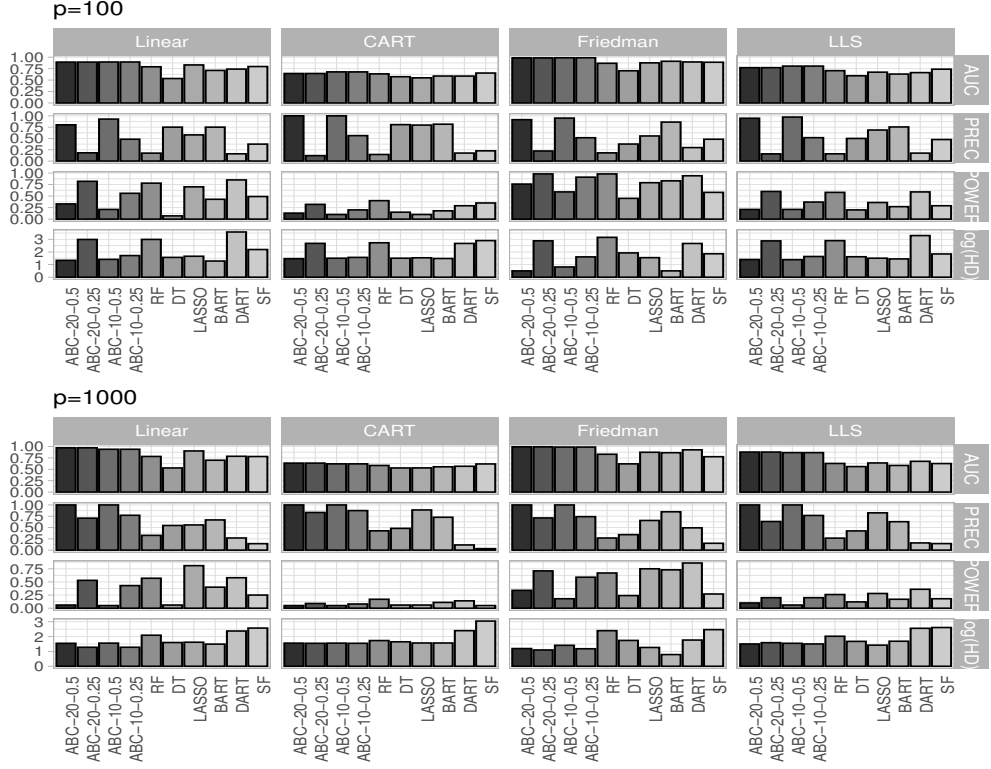


**Fig. 2.** Average variable selection performance under equicorrelation  $\rho_{ij} = 0.5$  over 20 simulations. Each panel corresponds to a different dimension  $p \in \{100, 1000\}$ . Each row reports a different statistic: AUC is the area under the ROC curve,  $\text{PREC} = 1 - \text{FDP} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ ,  $\text{POWER} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ ,  $\log(\text{HD}) = \log(\text{FP} + \text{FN})$ . ABC is run for  $T \in \{10, 20\}$  and cutoff  $\in \{0.5, 0.25\}$ . Each column indicates a different data generating process.

which have been split on inside a forest at least once on average. Alternative strategies based on truncating inclusion probabilities (Linero, 2018) using data-adaptive thresholds (Bleich et al., 2014) did not perform better, in general. For ABC, we report results for two selection thresholds 0.5 and 0.25. For Spike-and-Forest (SF), we report the median probability model.

The performance comparisons for variable selection are summarized in Figure 2 (equicorrelation  $\rho_{ij} = 0.5$ ) and Figure 3 (autocorrelation  $\rho_{ij} = 0.9^{|i-j|}$ ). These figures show that ABC has an advantage in terms of AUC, suggesting that ABC can rank variables more efficiently. While RF tend to have a higher power, they are plagued with false discoveries (i.e. smaller precision). ABC Bayesian Forests, on the other hand, are seen to yield fewer false discoveries (i.e. higher precision) relative to the other procedures. The ABC threshold 0.5 yields higher precision whereas 0.25 yields higher power.

While ABC Bayesian Forests were designed to explore the posterior distribution over models, it is natural to ask whether they also yield reasonable prediction. There are



**Fig. 3.** Average variable selection performance under autocorrelation  $\rho_{ij} = 0.9^{|i-j|}$  over 10 simulations. Each panel corresponds to a different dimension  $p \in \{100, 1000\}$ . Each row reports a different statistic: AUC is the area under the ROC curve,  $\text{PREC} = 1 - \text{FDP} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ ,  $\text{POWER} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ ,  $\log(\text{HD}) = \log(\text{FP} + \text{FN})$ . ABC is run for  $T \in \{10, 20\}$  and cutoff  $\in \{0.5, 0.25\}$ . Each column indicates a different data generating process.

various ways to perform prediction with our ABC method. One natural strategy is to save each draw  $f_{S,B}^m$  at the  $m^{\text{th}}$  ABC iteration when  $\epsilon_m < \epsilon$  and average out individual predictions obtained from these single draws. Alternatively, one could first select variables based on ABC Bayesian Forests and then run a separate BART method (using the default number of  $T = 200$  trees which is recommended for prediction) with the selected variables. Using both strategies, we report average out-of-sample mean squared prediction error, where the average is taken over 20 independent validation samples generated from the same data generating process (Table 1). We include both ABC predictions described above and denote them as ABC1 and ABC2, respectively, for the two different thresholds ( $c \in \{0.5, 0.25\}$ ) and for the two choices of the number of trees ( $T \in \{10, 20\}$ ).

The best method under each simulation setting is marked in bold. When the data becomes more non-linear (CART and LLS setups) and the correlation among variables gets stronger, ABC tends to outperform the other methods. DART, on the other hand, works better for more linear datasets. Note that our default ABC implementation internally

**Table 1.** Average out-of-sample mean squared prediction error over 20 independent validation datasets. ABC1 denotes predictions using ABC samples  $f_{S,B}^m$  and ABC2 uses ABC variable selection and runs BART ( $T = 200$ ) on the selected subset.  $T$  designates the number of trees and  $c$  is the selection threshold. The best performing method for each row is denoted in bold.

	ABC2 $T = 20$	ABC1 $T = 20, c = 0.5$	ABC1 $T = 20, c = 0.25$	ABC2 $T = 10$	ABC1 $T = 10, c = 0.5$	ABC1 $T = 10, c = 0.25$	RF	RLT	DT	BART	DART
<b>Equi-correlation <math>\rho_{ij} = 0.5</math> for <math>i \neq j</math></b>											
Linear											
$p = 100$	5.56	5.58	5.84	5.60	5.84	5.55	5.63	5.45	5.92	5.49	<b>5.40</b>
$p = 1000$	5.79	6.15	5.73	5.86	6.28	5.95	5.83	5.70	6.04	5.82	<b>5.62</b>
CART											
$p = 100$	34.21	34.63	37.19	<b>34.00</b>	36.10	35.81	34.21	34.64	34.61	35.48	35.57
$p = 1000$	32.00	34.27	35.72	<b>31.99</b>	33.93	33.17	32.30	32.40	33.08	33.77	34.04
Friedman											
$p = 100$	30.32	29.28	31.59	30.52	30.30	<b>29.03</b>	31.84	30.17	41.41	31.31	<b>29.03</b>
$p = 1000$	33.14	35.97	31.54	33.54	38.42	32.71	34.35	32.22	45.69	32.99	<b>29.42</b>
LLS											
$p = 100$	<b>26.23</b>	27.00	28.70	26.25	26.90	27.36	26.80	26.46	28.51	27.42	27.42
$p = 1000$	27.37	26.98	<b>26.94</b>	27.38	27.07	27.02	27.18	26.68	30.66	28.21	27.49
<b>Auto-correlation <math>\rho_{ij} = 0.9^{ i-j }</math></b>											
Linear											
$p = 100$	6.17	6.29	6.37	6.20	6.25	6.18	6.37	6.09	6.77	6.17	<b>5.91</b>
$p = 1000$	6.39	6.44	<b>6.00</b>	6.47	6.21	6.13	6.55	6.20	7.06	6.53	6.42
CART											
$p = 100$	33.80	37.72	37.28	33.83	36.78	36.61	<b>33.57</b>	34.40	35.05	35.61	35.81
$p = 1000$	31.57	33.55	37.21	<b>31.52</b>	33.52	37.43	31.63	31.88	32.22	33.11	33.43
Friedman											
$p = 100$	34.09	32.51	34.65	34.27	34.97	32.77	36.88	33.83	48.64	34.21	<b>30.36</b>
$p = 1000$	39.09	39.57	32.58	40.58	43.05	33.46	41.80	37.38	49.51	35.96	<b>30.81</b>
LLS											
$p = 100$	28.57	<b>27.94</b>	30.71	28.45	28.03	29.12	28.88	27.87	30.69	28.83	28.81
$p = 1000$	29.98	<b>28.25</b>	28.96	30.14	28.40	28.38	30.19	28.56	32.29	31.76	29.28

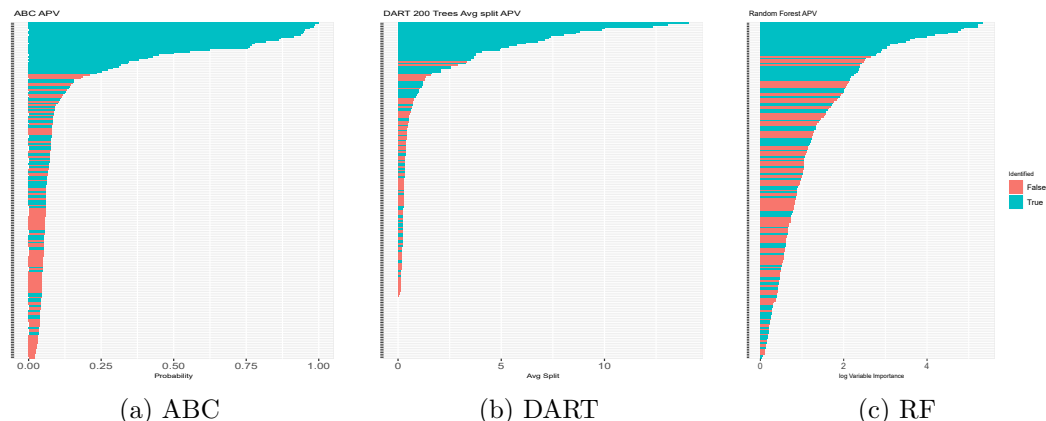
uses only a *small* number of  $B = 200$  burn-in iterations and a small number of trees. For prediction, it has been recommended that BART is deployed with a larger number of trees (Chipman et al., 2010). In addition, the ABC computation produces forest samples  $f_{S,B}^m$  which are from an *approximate* posterior. These two facts may affect resulting predictions which may not necessarily outperform BART (DART) across-the-board.

## 6. HIV Data

To further illustrate the usefulness of our approach, we consider a dataset described and analyzed in Rhee et al. (2006) and Barber and Candès (2015). The data consists of genotype and resistance measurements (log-decrease in susceptibility) for three drug classes, i.e. protease inhibitors (PIs), nucleoside reverse transcriptase inhibitors (NRTIs) and non-nucleoside reverse transcriptase inhibitors (NNRTIs). The data is publicly available from the Stanford HIV Drug Resistance Database.<sup>2</sup>

The goal of this analysis is to identify possible non-polymorphic mutation positions which result in a log-fold increase of lab-tested drug resistance. The design matrix  $X = (x_{ij})_{i,j=1}^{n,p}$  consists of binary indicators  $x_{ij} \in \{0, 1\}$  for whether or not the  $j^{\text{th}}$  mutation occurred in the  $i^{\text{th}}$  sample. As in Barber and Candès (2015), only mutations that appear at least 3 times are taken into consideration. One appealing feature of this dataset is

<sup>2</sup> [https://hivdb.stanford.edu/pages/published\\_analysis/genophenoPNAS2006/](https://hivdb.stanford.edu/pages/published_analysis/genophenoPNAS2006/)

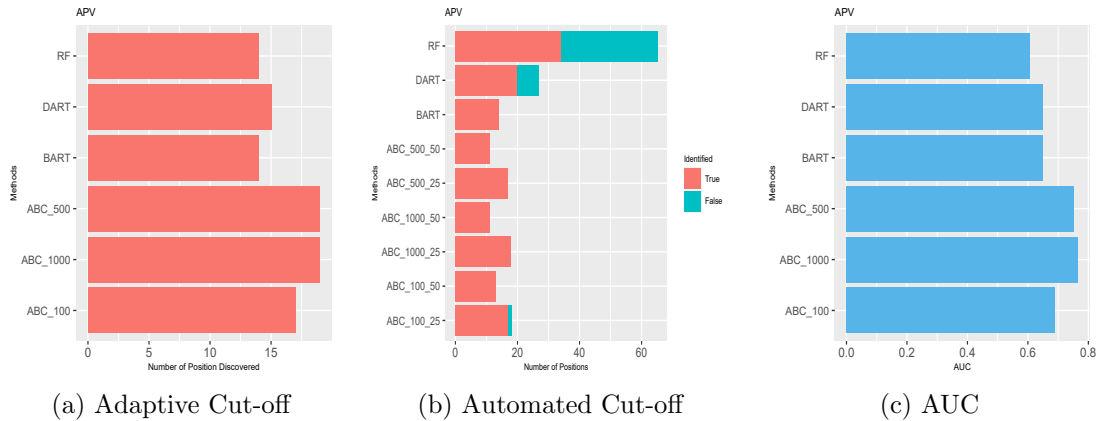


**Fig. 4.** A barplot of ordered importance measures (inclusion probabilities for ABC, importance measures for DART and RF) for each of the  $p = 201$  mutations for the drug APV, where blue represents mutations found in Rhee et al. (2005). (a) Inclusion probabilities are computed using the top 1000 out of  $M = 10\,000$  ABC samples; (b) Average split of DART with 20 000 MCMC iterations; (c) log variable importance of Random Forest with 500 trees.

the availability of a proxy to the ‘ground truth’. Indeed, in an independent experimental study, Rhee et al. (2005) identified mutations that are present at a significantly higher frequency in patients who have been treated with each drug. Similarly as Barber and Candès (2015), we treat this experimental data as an approximation to the truth for comparisons and for validation of our findings.

We run ABC with  $M = 10\,000$  iterations, where each internal BART sample is obtained after 200 burnin iterations with 20 trees. The top 1000 ABC samples with the smallest  $\epsilon_m$  are kept and used to compute inclusion probabilities for each mutation. For illustration, we visualize results for one of the PI drugs (APV) and report the results for all the drugs in the Supplemental Material (Section S.5). The inclusion probabilities have been ordered and plotted in Figure 4, where the mutations experimentally validated by Rhee et al. (2005) (a proxy for true signals) are denoted in blue and the rest is in red. For comparisons, we also included the importance measure (the average number of splits on each variable) from DART run with 20 000 MCMC iterations and  $T = 200$  trees as well as the importance measure (on a log scale) from Random Forests (RF) run with 500 trees.

Figure 4 reveals that ABC Bayesian Forests have a strong separation power, where experimentally validated mutations generally have a higher inclusion probability. Compared to DART and RF, ABC clearly stands out as being more effective in weeding out ‘noise’. We gauge the strength of the signal/noise separation using several descriptive statistics. In these comparisons, we also consider plain BART method (using  $T = 20$  trees and 20 000 MCMC iterations) and ABC using the top 100 and 500 samples with the smallest tolerance level  $\epsilon_m$ . Since the selection of the cut-off point is not obvious for BART and RF, we first select variables based on an adaptive cut-off point so that there are no false discoveries (i.e. the cut-off is the largest importance weight of a *not*



**Fig. 5.** (a) The number of true discoveries using an adaptive cut-off; (b) The number of true (red) and false (blue) discoveries using an automated cut-off; (c) The AUC of each method.

experimentally validated mutation). From the plot of the number of ‘True’ locations selected (displayed in Figure 5(a)) we can see that all three ABC implementations find more signal variables. Next, we choose the cut-off point in an automated way, where ABC importance probabilities are truncated at 0.5 and 0.25, BART and DART measures are truncated at one (i.e. the variable has been used on average at least once), and RF select variables using recursive feature elimination as explain in the previous section. Similarly to Barber and Candès (2015), we report the number of ‘True’ locations and ‘False’ locations (Figure 5(b)). RF selection is plagued with false discoveries and DART is not free from false identifications either. The ABC selection cutoff 0.5 results in a more conservative selection, where lowering the cutoff point to 0.25 yields more discoveries. Finally, from the plot of the AUC values for all considered methods (Figure 5(c)), we conclude that ABC is better at separating the experimentally validated mutations from the rest even using a very few filtered ABC samples.

## 7. Discussion

This paper makes advancements at two fronts. One is the proposal of ABC Bayesian Forests for variable selection based on a new idea of data splitting, where a fraction of data is first used for ABC proposal and the rest for ABC rejection. This new strategy increases ABC acceptance rate. We have shown that ABC Bayesian Forests are highly competitive with (and often better than) other tree-based variable selection procedures. The second development is theoretical and concerns consistency for variable and regularity selection. Continuing the theoretical investigation of BART by Ročková and van der Pas (2017), we proposed new complexity priors which jointly penalize model dimensionality and tree size. We have shown joint consistency for variable *and* regularity selection when the level of smoothness is unknown and no greater than 1. Our results are the first model selection consistency results for BART priors.

Our ABC sampling routine has the potential to be extended in various ways. Sampling from  $\pi(f_{\mathcal{E},B}, \sigma^2 | \mathbf{Y}_{\mathcal{I}_m}^{obs}, \mathcal{S}_m)$  in ABC Bayesian Forests is one way of distilling  $\mathbf{Y}_{\mathcal{I}_m}^{obs}$

to propose a candidate ensemble  $f_{\mathcal{E},\mathbf{B}}^m$ . We noticed that the ABC acceptance rate can be further improved by replacing a randomly sampled tree with a fitted tree. Indeed, instead of drawing from  $\pi(f_{\mathcal{E},\mathbf{B}}, \sigma^2 | \mathbf{Y}_{\mathcal{I}_m}^{obs}, \mathcal{S})$ , one can fit a tree  $\hat{f}_{\mathcal{T},\mathbf{B}}^m$  to  $\mathbf{Y}_{\mathcal{I}_m}^{obs}$  using recursive partitioning algorithms (such as the `rpart` R package of [Therneau and Atkinson \(2018\)](#) or with BART (by taking the posterior mean estimate  $\hat{f}_{\mathcal{E},\mathbf{B}}^m = \mathbb{E}[f_{\mathcal{E},\mathbf{B}} | \mathbf{Y}_{\mathcal{I}_m}^{obs}, \mathcal{S}]$ ). This variant, further referred to as ABC Forest Fit, is indirectly linked to other model-selection methods based on resampling.

[Felsenstein \(1985\)](#) proposed a “first-order bootstrap” to assess confidence of an estimated tree phylogeny. The idea was to construct a tree from each bootstrap sample and record the proportion of bootstrap trees that have a feature of interest (for us, this would be variables used for splits). [Efron and Tibshirani \(1998\)](#) embedded this approach within a parametric bootstrap framework, linking the bootstrap confidence level to both frequentist  $p$ -values and Bayesian a posteriori model probabilities. The authors proposed a second-order extension by reweighting the first-order resamples according to a simple importance sampling scheme. This second-order variant performs frequentist calibration of the a-posteriori probabilities and amounts to performing Bayesian analysis with Welch-Peers uninformative priors. [Efron \(2012\)](#) further develops the connection between parametric Bootstrap and posterior sampling through reweighting in exponential family models. Using non-parametric bootstrap ideas, [Newton and Raftery \(1994\)](#) introduce the weighted likelihood bootstrap (WLB) to sample from approximate posterior distributions. The WLB samples are obtained by maximum reweighted likelihood estimation with random weights. Such posterior sampling can be beneficial when, for instance, maximization is easier than Gibbs sampling from conditionals. In a similar spirit, our ABC Forest Fit variant would perform optimization (instead of sampling) on a random subset of the dataset to obtain a candidate tree/ensemble.

It is worth pointing out that  $\hat{f}_{\mathcal{E},\mathbf{B}}^m$  does not necessarily have to be a tree/forest. We suggest trees because they are easily trainable and produce stable results using traditional software packages. In principle, however, this method could be deployed in tandem with other non-parametric methods, such as deep learning, to perform variable selection.

## Acknowledgments

This work was supported by the James S. Kemper Foundation Faculty Research Fund at the University of Chicago Booth School of Business.

## References

- Barber, R. F. and Candès, E. J. (2015) Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, **43**, 2055–2085.
- Barber, S., Voss, J. and Webster, M. (2015) The rate of convergence for approximate Bayesian computation. *Electronic Journal of Statistics*, **9**, 80–105.
- Barbieri, M. M. and Berger, J. O. (2004) Optimal predictive model selection. *The Annals of Statistics*, **32**, 870–897.



- Berger, J. O. and Pericchi, L. R. (1996) The intrinsic Bayes factor for linear models. *Bayesian statistics*, **5**, 25–44.
- (2004) Training samples in objective Bayesian model selection. *The Annals of Statistics*, **32**, 841–869.
- Bleich, J., Kapelner, A., George, E. I. and Jensen, S. T. (2014) Variable selection for BART: An application to gene regulation. *The Annals of Applied Statistics*, 1750–1781.
- Breiman, L. (2001) Random forests. *Machine Learning*, **45**, 5–32.
- Breiman, L. and Cutler, A. (2013) Online manual for random forests. URL: [www.stat.berkeley.edu/~breiman/RandomForests/cc\\_home.html](http://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.html).
- Breiman, L., Friedman, J., Olshen, R. and Stone, C. J. (1984) *Classification and regression trees*. New York: Chapman and Hall.
- Candes, E., Fan, Y., Janson, L. and Lv, J. (2018) Panning for gold: ‘model-x’ knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- Carbonetto, P. and Stephens, M. (2012) Scalable variational inference for bayesian variable selection in regression, and its accuracy in genetic association studies. *Bayesian Analysis*, **7**, 73–108.
- Castillo, I. and Misner, R. (2018) Empirical Bayes analysis of spike and slab posterior distributions. *arXiv preprint arXiv:1801.01696*.
- Castillo, I., Schmidt-Hieber, J. and Van der Vaart, A. (2015) Bayesian linear regression with sparse priors. *The Annals of Statistics*, **43**, 1986–2018.
- Castillo, I. and van der Vaart, A. (2012) Needles and straw in a haystack: Posterior concentration for possibly sparse sequences. *The Annals of Statistics*, **40**, 2069–2101.
- Chipman, H., George, E. I. and McCulloch, R. E. (2001) The practical implementation of Bayesian model selection. In *Model Selection*, 65–116. Institute of Mathematical Statistics.
- Chipman, H. A., George, E. I. and McCulloch, R. E. (1998) Bayesian CART model search. *Journal of the American Statistical Association*, **93**, 935–948.
- (2010) BART: Bayesian additive regression trees. *The Annals of Applied Statistics*, **4**, 266–298.
- Comminges, L. and Dalalyan, A. S. (2012) Tight conditions for consistency of variable selection in the context of high dimensionality. *The Annals of Statistics*, **40**, 2667–2696.
- Csillery, K., Blum, M. G., Gaggiotti, O. E. and Francois, O. (2010) Approximate Bayesian computation (ABC) in practice. *Trends in ecology & evolution*, **25**, 410–418.

- Denison, D. G., Mallick, B. K. and Smith, A. F. (1998) A Bayesian CART algorithm. *Biometrika*, **85**, 363–377.
- Efron, B. (2012) Bayesian inference and the parametric bootstrap. *The Annals of Applied Statistics*, **6**, 1971.
- Efron, B. and Tibshirani, R. (1998) The problem of regions. *The Annals of Statistics*, 1687–1718.
- Fan, J. and Li, R. (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, **96**, 1348–1360.
- Felsenstein, J. (1985) Confidence limits on phylogenies: an approach using the bootstrap. *Evolution*, **39**, 783–791.
- Frazier, D. T., Martin, G. M., Robert, C. P. and Rousseau, J. (2018) Asymptotic properties of approximate Bayesian computation. *Biometrika*, **105**, 593–607.
- Frazier, D. T., Robert, C. P. and Rousseau, J. (2020) Model misspecification in approximate Bayesian computation: consequences and diagnostics. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- Friedman, J., Hastie, T. and Tibshirani, R. (2010) Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, **33**, 1.
- Friedman, J. H. (1991) Multivariate adaptive regression splines. *The Annals of Statistics*, **19**, 1–67.
- George, E. I. and McCulloch, R. E. (1993) Variable selection via Gibbs sampling. *Journal of the American Statistical Association*, **88**, 881–889.
- Ghosal, S., Lember, J. and Van Der Vaart, A. (2008) Nonparametric bayesian model selection and averaging. *Electronic Journal of Statistics*, **2**, 63–89.
- Ghosal, S. and van der Vaart, A. (2007) Convergence rates of posterior distributions for noniid observations. *The Annals of Statistics*, **35**, 192–223.
- Ghosh, J. K. and Samanta, T. (2002) Nonsubjective bayes testing?an overview. *Journal of statistical planning and inference*, **103**, 205–223.
- Good, I. J. (1950) Probability and the weighing of evidence.
- Gramacy, R. and Lee, H. (2008) Bayesian treed gaussian process models with an application to computer modeling. *Journal of the American Statistical Association*, **103**, 1119–11303.
- Grelaud, A., Robert, C. P. and Marin, J.-M. (2009) ABC methods for model choice in Gibbs random fields. *Comptes Rendus Mathematique*, **347**, 205–210.
- Guan, Y. and Stephens, M. (2011) Bayesian variable selection regression for genome-wide association studies and other large-scale problems. *The Annals of Applied Statistics*, **5**, 1780–1815.

- Hill, J. (2011) Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, **20**, 217–240.
- Ishwaran, H. (2007) Variable importance in binary regression trees and forests. *Electronic Journal of Statistics*, **1**, 519–537.
- Jiang, B., Wu, T., Zheng, C. and Wong, W. (2017) Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica*, **27**, 1595–1618.
- Johnson, V. E. and Rossell, D. (2012) Bayesian model selection in high-dimensional settings. *Journal of the American Statistical Association*, **107**, 649–660.
- Kazemitabar, J., Amini, A., Bloniarz, A. and Talwalkar, A. S. (2017) Variable importance using decision trees. In *Advances in Neural Information Processing Systems*, 425–434.
- Lafferty, J. and Wasserman, L. (2001) Iterative Markov Chain Monte Carlo computation of reference priors and minimax risk. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, 293–300. Morgan Kaufmann Publishers Inc.
- (2008) RODEO: sparse, greedy nonparametric regression. *The Annals of Statistics*, 28–63.
- Lember, J. and van der Vaart, A. (2007) On universal Bayesian adaptation. *Statistics & Decisions*, **25**, 127–152.
- Lempers, F. B. (1971) Posterior probabilities of alternative linear models.
- Li, W. and Fearnhead, P. (2018) Convergence of regression-adjusted approximate Bayesian computation. *Biometrika*, **105**, 301–318.
- Liang, F., Li, Q. and Zhou, L. (2018) Bayesian neural networks for selection of drug sensitive genes. *Journal of the American Statistical Association*, **113**, 955–972.
- Liaw, A. and Wiener, M. (2002) Classification and regression by randomForest. *R news*, **2**, 18–22.
- Lin, Y. and Zhang, H. H. (2006) Component selection and smoothing in multivariate nonparametric regression. *The Annals of Statistics*, **34**, 2272–2297.
- Linero, A. R. (2018) Bayesian regression trees for high-dimensional prediction and variable selection. *Journal of the American Statistical Association*, 1–11.
- Marin, J.-M., Pudlo, P., Robert, C. P. and Ryder, R. J. (2012) Approximate Bayesian computational methods. *Statistics and Computing*, **22**, 1167–1180.
- Martin, J. S., Jasra, A., Singh, S. S., Whiteley, N., Del Moral, P. and McCoy, E. (2014) Approximate Bayesian computation for smoothing. *Stochastic Analysis and Applications*, **32**, 397–420.

- McCulloch, R., Sparapani, R., Gramacy, R., Spanbauer, C. and Pratola, M. (2018) *BART: Bayesian Additive Regression Trees*. URL: <https://CRAN.R-project.org/package=BART>. R package version 1.6.
- Moreno, E., Girón, J. and Casella, G. (2015) Posterior model consistency in variable selection as the model dimension grows. *Statistical Science*, **30**, 228–241.
- Narisetty, N. N. and He, X. (2014) Bayesian variable selection with shrinking and diffusing priors. *The Annals of Statistics*, **42**, 789–817.
- Newton, M. A. and Raftery, A. E. (1994) Approximate Bayesian inference with the weighted likelihood bootstrap. *Journal of the Royal Statistical Society. Series B (Methodological)*, 3–48.
- O’Hagan, A. (1995) Fractional bayes factors for model comparison. *Journal of the Royal Statistical Society: Series B (Methodological)*, **57**, 99–118.
- Pérez, J. M. and Berger, J. O. (2002) Expected-posterior prior distributions for model selection. *Biometrika*, **89**, 491–512.
- Plagnol, V. and Tavaré, S. (2004) Approximate Bayesian Computation and MCMC. In *Monte Carlo and Quasi-Monte Carlo Methods 2002*, 99–113. Springer.
- Pratola, M. T. (2016) Efficient Metropolis–Hastings proposal mechanisms for Bayesian regression tree models. *Bayesian Analysis*, **11**, 885–911.
- Pudlo, P., Marin, J.-M., Estoup, A., Cornuet, J.-M., Gautier, M. and Robert, C. P. (2015) Reliable ABC model choice via random forests. *Bioinformatics*, **32**, 859–866.
- Radchenko, P. and James, G. M. (2010) Variable selection using adaptive nonlinear interaction structures in high dimensions. *Journal of the American Statistical Association*, **105**, 1541–1553.
- Ravikumar, P., Lafferty, J., Liu, H. and Wasserman, L. (2009) Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **71**, 1009–1030.
- Rhee, S.-Y., Fessel, W. J., Zolopa, A. R., Hurley, L., Liu, T., Taylor, J., Nguyen, D. P., Slome, S., Klein, D. and Horberg, M. (2005) Hiv-1 protease and reverse-transcriptase mutations: correlations with antiretroviral therapy in subtype b isolates and implications for drug-resistance surveillance. *The Journal of infectious diseases*, **192**, 456–465.
- Rhee, S.-Y., Taylor, J., Wadhera, G., Ben-Hur, A., Brutlag, D. L. and Shafer, R. W. (2006) Genotypic predictors of human immunodeficiency virus type 1 drug resistance. *Proceedings of the National Academy of Sciences*, **103**, 17355–17360.
- Robert, C. P., Cornuet, J.-M., Marin, J.-M. and Pillai, N. S. (2011) Lack of confidence in approximate Bayesian computation model choice. *Proceedings of the National Academy of Sciences*, **108**, 15112–15117.
- Ročková, V. and George, E. I. (2014) EMVS: The EM approach to Bayesian variable selection. *Journal of the American Statistical Association*, **109**, 828–846.

- (2018) The Spike-and-Slab LASSO. *Journal of the American Statistical Association*, **113**, 431–444.
- Ročková, V. and van der Pas, S. (2017) Posterior concentration for Bayesian regression trees and their ensembles. *The Annals of Statistics (in revision)*.
- Ročková, V. (2017) Particle EM for variable selection. *Journal of the American Statistical Association*, 1–30.
- Ročková, V. and Saha, E. (2019) On theory for BART. In *Artificial Intelligence and Statistics*.
- Savitsky, T., Vannucci, M. and Sha, N. (2011) Variable selection for nonparametric Gaussian process priors: Models and computational strategies. *Statistical Science*, **26**, 130.
- Scheipl, F. (2011) spikeslabgam: Bayesian variable selection, model choice and regularization for generalized additive mixed models in r. *arXiv preprint arXiv:1105.5253*.
- Scott, J. G. and Berger, J. O. (2010) Bayes and empirical-Bayes multiplicity adjustment in the variable-selection problem. *The Annals of Statistics*, 2587–2619.
- Sunnaaker, M., Busetto, A. G., Numminen, E., Corander, J., Foll, M. and Dessimoz, C. (2013) Approximate bayesian computation. *PLoS computational biology*, **9**, e1002803.
- Taddy, M., Gramacy, R. and Polson, N. (2011a) Dynamic trees for learning and design. *Journal of the American Statistical Association*, **106**, 409–123.
- Taddy, M. A., Gramacy, R. B. and Polson, N. G. (2011b) Dynamic trees for learning and design. *Journal of the American Statistical Association*, **106**, 109–123.
- Tavaré, S., Balding, D. J., Griffiths, R. C. and Donnelly, P. (1997) Inferring coalescence times from DNA sequence data. *Genetics*, **145**, 505–518.
- Therneau, T. and Atkinson, B. (2018) *rpart: Recursive Partitioning and Regression Trees*. URL: <https://CRAN.R-project.org/package=rpart>. R package version 4.1-13.
- Tibshirani, R. (1996) Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 267–288.
- Turlach, B. A. (2004) Discussion on least angle regression. *The Annals of Statistics*, **32**, 481–490.
- Van De Geer, S. A. and Bühlmann, P. (2009) On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, **3**, 1360–1392.
- Yang, Y. and Pati, D. (2017) Bayesian model selection consistency and oracle inequality with intractable marginal likelihood. *arXiv preprint arXiv:1701.00311*.
- Zhao, P. and Yu, B. (2006) On model selection consistency of LASSO. *Journal of Machine learning research*, **7**, 2541–2563.

Zhu, R., Zeng, D. and Kosorok, M. R. (2015) Reinforcement learning trees. *Journal of the American Statistical Association*, **110**, 1770–1784.

# Supplemental Materials

## S.1. Theory

### S.1.1. Proof of Theorem 4.1

We first review some notation used throughout this section and adapted from Ročková and van der Pas (2017). Recall that  $\Pi_{\mathcal{S}}(\cdot)$  denotes the conditional distribution given the model  $\mathcal{S}$ . Next,  $\mathcal{F}_{\mathcal{S}}(K)$  denotes a set of all step functions  $f_{\mathcal{T},\beta}(\cdot)$  with  $K$  steps that split on covariates  $\mathcal{S}$  and  $\|f_{\mathcal{T},\beta}\|_{\infty} \leq B$ . A tree partition is called valid when each tree splits on observed values  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and has nonempty cells. We denote with  $\mathcal{V}_{\mathcal{S}}^K$  all valid trees obtained by splitting  $K - 1$  times along coordinates inside  $\mathcal{S}$ . The number of such valid trees is denoted with  $\Delta(\mathcal{V}_{\mathcal{S}}^K)$ . For a valid tree partition  $\mathcal{T} \in \mathcal{V}_{\mathcal{S}}^K$ , we denote with  $\mathcal{F}(\mathcal{T}) \subset \mathcal{F}_{\mathcal{S}}(K)$  all step functions supported on  $\mathcal{T}$ . We prove Theorem 4.1 by verifying conditions B1-B4 in Theorem 4 of Yang and Pati (2017) (further referred to as YP17). We build on tools developed in Ročková and van der Pas (2017) (further referred to as RP17).

#### S.1.1.1. Prior Concentration Condition

The first condition pertains to prior concentration and consists of two parts: (a) the model prior mass condition and (b) the prior concentration condition in the parameter space under the true model. Namely, we want to show that

$$\pi(\mathcal{S}_0) \geq e^{-n\varepsilon_{n,\mathcal{S}_0}^2} \quad (26)$$

and

$$\Pi_{\mathcal{S}_0}(f_{\mathcal{T},\beta} \in \mathcal{F}_{\mathcal{S}_0}(K) : \|f_{\mathcal{T},\beta} - f_0\|_n \leq \varepsilon_{n,\mathcal{S}_0}) \geq e^{-dn\varepsilon_{n,\mathcal{S}_0}^2} \quad (27)$$

for some  $d > 2$ . The prior concentration (26) follows directly from the definition of model weights (13) for  $C \leq C_{\varepsilon}^2$  under our assumption  $q_0 \log p < n^{q_0/(2\alpha+q_0)}$ .

Regarding (27), a variant of this condition is verified in Section 8.2 of RP17 assuming that  $K$  is random with a prior. It follows from their proof, however, that (27) holds if we fix  $K$  at  $K_{\mathcal{S}_0} = \lfloor C_K/C_{\varepsilon}^2 n \varepsilon_{n,\mathcal{S}_0}^2 / \log n \rfloor = 2^{q_0 s}$  for some  $s \in \mathbb{N}$ . The proof consists of (a) constructing a single approximating tree (i.e. the  $k$ - $d$  tree with  $s = (\log_2 K_{\mathcal{S}_0})/q_0$  cycles of splits on each coordinate in  $\mathcal{S}_0$ ) and showing that it has enough prior support. This tree exists under the assumption that the design is  $\mathcal{S}_0$ -regular. From (8.5) of RP17, such tree approximates  $f_0$  with an error bounded by a constant multiple of  $\varepsilon_{n,\mathcal{S}_0}$ . The verification of (27) then follows directly from RP17.

#### S.1.1.2. Entropy Condition

The second condition (B4 in the notation of YP17) entails controlling the complexity of over/underfitting models. In the sequel, we focus only on models with up to  $q_n$  covariates, where  $q_n = C_q \lceil n \varepsilon_{n,\mathcal{S}_0}^2 / \log p \rceil$ . This restriction is justified by the following lemma.



LEMMA S.1.1. Denote with  $q_n = C_q \lceil n \varepsilon_{n, \mathcal{S}_0}^2 / \log p \rceil$ . Under the assumptions of Theorem 4.1, we have

$$\Pi(q \geq q_n \mid \mathbf{Y}^{(n)}) \rightarrow 0 \quad (28)$$

in  $\mathbb{P}_{f_0}^{(n)}$ -probability as  $n \rightarrow \infty$ .

PROOF. First, we show that  $\Pi(q \geq q_n) e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \rightarrow 0$ , where  $d > 2$  is as in (27). We can write

$$\begin{aligned} \Pi(q > q_n) e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} &\lesssim e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \sum_{k=q_n}^p \binom{p}{k} e^{-C \times \max\{n^{k/(2\alpha+k)} \log n, k \log p\}} \\ &\leq e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2 - (C-2)q_n \log p} = e^{-n \varepsilon_{n, \mathcal{S}_0}^2 [(C-2)C_q - (d+2)]}. \end{aligned}$$

The right hand side above goes to zero when  $(C-2)C_q - (d+2) > 0$ . This can be satisfied with  $C > 2$  and  $C_q$  large enough. This fact, together with prior mass conditions (27) and (26), yields (28) according to Lemma 1 of Ghosal and van der Vaart (2007).

Lemma S.1.1 essentially states that the posterior will not reward models whose dimensionality is larger than (or equal to)  $q_n$ . In our following considerations, we thus condition only models with less than  $q_n$  variables.

We now verify that the complexity of overfitting models  $\mathcal{S} \supset \mathcal{S}_0$  is not too large in the sense that their global metric entropy satisfies

$$\log N(\varepsilon_{n, \mathcal{S}}; \mathcal{F}_{\mathcal{S}}(K_{\mathcal{S}}); \|\cdot\|_n) \leq n \varepsilon_{n, \mathcal{S}}^2. \quad (29)$$

First, we note that for two tree step functions  $f_{\mathcal{T}, \beta_1} \in \mathcal{F}(\mathcal{T})$  and  $f_{\mathcal{T}, \beta_2} \in \mathcal{F}(\mathcal{T})$  that have the same partition  $\mathcal{T} \in \mathcal{V}_{\mathcal{S}}^{K_{\mathcal{S}}}$  and different step heights  $\beta_1 \in \mathbb{R}^{K_{\mathcal{S}}}$  and  $\beta_2 \in \mathbb{R}^{K_{\mathcal{S}}}$ , we have  $\{\|f_{\mathcal{T}, \beta_1} - f_{\mathcal{T}, \beta_2}\|_n \leq \varepsilon_{n, \mathcal{S}}\} \supset \{\|\beta_1 - \beta_2\|_2 \leq \varepsilon_{n, \mathcal{S}}\}$ . Furthermore, noting that  $\mathcal{F}(\mathcal{T}) = \{f_{\mathcal{T}, \beta} : \|f_{\mathcal{T}, \beta}\|_{\infty} \leq B\} \subset \{\beta \in \mathbb{R}^{K_{\mathcal{S}}} : \|\beta\|_2 \leq B \sqrt{n}\}$  we can write

$$N(\varepsilon_{n, \mathcal{S}}; \mathcal{F}(\mathcal{T}); \|\cdot\|_n) \leq \left( \frac{3B \sqrt{n}}{\varepsilon_{n, \mathcal{S}}} \right)^{K_{\mathcal{S}}} \leq \left( 3B n^{3/2} / C_{\varepsilon} \right)^{K_{\mathcal{S}}},$$

where we used the standard  $\varepsilon_{n, \mathcal{S}}$  covering number of a  $K_{\mathcal{S}}$ -Euclidean ball of a radius  $B \sqrt{n}$  and the fact that  $1/\varepsilon_{n, \mathcal{S}} \leq 1/C_{\varepsilon} \times n^{\alpha/(2\alpha+|\mathcal{S}|)} \leq 1/C_{\varepsilon} \times n$ . Then we can write

$$N(\varepsilon_{n, \mathcal{S}}; \mathcal{F}_{\mathcal{S}}(K_{\mathcal{S}}); \|\cdot\|_n) \leq \Delta(\mathcal{V}_{\mathcal{S}}^{K_{\mathcal{S}}}) \left( 3B n^{3/2} / C_{\varepsilon} \right)^{K_{\mathcal{S}}}.$$

Using Lemma 3.1 of Rockova and van der Pas (2017), we have  $\Delta(\mathcal{V}_{\mathcal{S}}^{K_{\mathcal{S}}}) \leq (K_{\mathcal{S}} n |\mathcal{S}|)^{K_{\mathcal{S}}}$ .

The overall log-covering number is then upper-bounded with (since  $|\mathcal{S}| \leq q_n \leq n$ )

$$K_{\mathcal{S}} \log \left( 3B n^3 n^{3/2} \right) \lesssim K_{\mathcal{S}} \log n \propto n \varepsilon_{n, \mathcal{S}}^2. \quad (30)$$

This verifies the model complexity condition for overfitting models. Next, we need to verify (29) with  $\varepsilon_{n, \mathcal{S}}$  replaced by  $\tilde{\varepsilon}_n$  for ‘‘underfitting’’ models  $\mathcal{S} \in \mathbf{\Gamma}_{\mathcal{S} \not\supset \mathcal{S}_0}$  where  $|\mathcal{S}| \leq q_n$ .

This follows from the same arguments as above and the fact that  $\varepsilon_{n,S} \leq \tilde{\varepsilon}_n$ . Finally, the last requirement in Assumption B4 of YP17 is verifying that

$$\sum_{\mathcal{S} \not\supset \mathcal{S}_0: |\mathcal{S}| \leq q_n} e^{-C_2 n \tilde{\varepsilon}_n^2} + \sum_{\mathcal{S} \supset \mathcal{S}_0: |\mathcal{S}| \leq q_n} e^{-C_2 n \varepsilon_{n,S}^2} \leq 1 \quad (31)$$

for some large constant  $C_2 > 0$ . Since  $\tilde{\varepsilon}_n \geq \varepsilon_{n,S} > \varepsilon_{n,\mathcal{S}_0}$  for any  $\mathcal{S} \supset \mathcal{S}_0$  such that  $|\mathcal{S}| \leq q_n$ , we can upper-bound the left-hand side above with

$$\sum_{q=0}^{q_n} \sum_{\mathcal{S}: |\mathcal{S}|=q} e^{-C_2 n \varepsilon_{n,\mathcal{S}_0}^2} \leq e^{-C_2 n \varepsilon_{n,\mathcal{S}_0}^2} \sum_{q=0}^{q_n} \binom{p}{q} \leq \left( \frac{2ep}{q_n} \right)^{q_n+1} e^{-C_2 n \varepsilon_{n,\mathcal{S}_0}^2}$$

From our definition of  $q_n$ , we have  $q_n \log p \asymp n \varepsilon_{n,\mathcal{S}_0}^2$  and (31) will be satisfied for a large enough  $C_2$ .

### S.1.1.3. Prior Anticoncentration Condition

Lastly, as one of the sufficient conditions for model selection consistency, we need to verify

$$\sum_{\mathcal{S} \supset \mathcal{S}_0: |\mathcal{S}| \leq q_n} \pi(\mathcal{S}) \Pi_{\mathcal{S}}(f_{\mathcal{T},\beta} \in \mathcal{F}_{\mathcal{S}}(K_{\mathcal{S}}) : \|f_0 - f_{\mathcal{T},\beta}\|_n \leq M \varepsilon_{n,\mathcal{S}}) \leq e^{-H n \varepsilon_{n,\mathcal{S}_0}^2} \quad (32)$$

for some  $H > 0$ . Alternatively, YP17 introduce the so-called ‘‘anti-concentration condition’’  $\Pi_{\mathcal{S}}(f_{\mathcal{T},\beta} \in \mathcal{F}_{\mathcal{S}}(K_{\mathcal{S}}) : \|f_0 - f_{\mathcal{T},\beta}\|_n \leq M \varepsilon_{n,\mathcal{S}}) \leq e^{-H n \varepsilon_{n,\mathcal{S}_0}^2}$  for overfitting models  $\mathcal{S} \supset \mathcal{S}_0$  where  $\varepsilon_{n,\mathcal{S}} \geq \varepsilon_{n,\mathcal{S}_0}$ . This condition is needed to show that the posterior probability of more complex models that contain the truth goes to zero.

It turns out that this condition can be avoided with our choice of model weights  $\pi(\mathcal{S})$  (Ghosal et al., 2008). We can verify (32) directly (without the anticoncentration condition) by upper-bounding the left hand side of (32) with

$$\sum_{\mathcal{S} \supset \mathcal{S}_0: |\mathcal{S}| \leq q_n} \pi(\mathcal{S}) \leq \sum_{\mathcal{S} \supset \mathcal{S}_0: |\mathcal{S}| \leq q_n} e^{-C n \varepsilon_{n,S}^2} \leq e^{-C n \varepsilon_{n,\mathcal{S}_0}^2} \left( \frac{2ep}{q_n} \right)^{q_n+1}. \quad (33)$$

Since  $q_n \log p \asymp n \varepsilon_{n,\mathcal{S}_0}^2$ , (32) holds for  $H < C - 1$ .

### S.1.1.4. Identifiability

Under the identifiability and irrepresentability assumptions (4.1) and (4.2), it turns out that we cannot approximate  $f_0$  well enough with models that miss at least one covariate. This property is summarized in the following Lemma, which is a variant of Proposition 1 of YP17.

LEMMA S.1.2. *For  $f_0 \in \mathcal{H}_p^\alpha \cap \mathcal{C}(\mathcal{S}_0)$ , assume that  $\mathcal{S}_0$  is  $(f_0, \varepsilon)$ -identifiable and that  $\varepsilon$ -irrepresentability holds. Then*

$$\inf_{\mathcal{S} \not\supset \mathcal{S}_0} \inf_{f_{\mathcal{T},\beta} \in \mathcal{F}_{\mathcal{S}}} \|f_0 - f_{\mathcal{T},\beta}\|_n > M \varepsilon.$$

PROOF. We decompose  $\mathcal{S} \not\supseteq \mathcal{S}_0$  into true positives and false positives, i.e.  $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$ , where  $\mathcal{S}_1 \subset \mathcal{S}_0$  and  $\mathcal{S}_2 \cap \mathcal{S}_0 = \emptyset$ . We denote with  $f^{\mathcal{S}}$  the projection of  $f_0$  onto  $\mathcal{F}_{\mathcal{S}}$ , omitting the subscripts  $\widehat{\mathcal{T}}$  and  $\widehat{\beta}$ . With a slight abuse of notation we denote  $\mathbb{E}(f, g) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)g(\mathbf{x}_i)$ . Then we can write

$$\|f_0 - \widehat{f}^{\mathcal{S}}\|_n^2 = \|f_0 - \widehat{f}^{\mathcal{S}_1} + \widehat{f}^{\mathcal{S}_1} - \widehat{f}^{\mathcal{S}}\|_n^2 > \|f_0 - \widehat{f}^{\mathcal{S}_1}\|_n^2 - 2|\mathbb{E}[(f_0 - \widehat{f}^{\mathcal{S}_1})(\widehat{f}^{\mathcal{S}} - \widehat{f}^{\mathcal{S}_1})]|,$$

where  $\mathbb{E}[(f_0 - \widehat{f}^{\mathcal{S}_1})(\widehat{f}^{\mathcal{S}} - \widehat{f}^{\mathcal{S}_1})]$  equals  $\rho_n^{\mathcal{S}}$  defined in (18). We note that  $\delta_n^{\mathcal{S}_1}$  is monotone increasing in the number of false non-discoveries  $|\mathcal{S}_0 \setminus \mathcal{S}_1|$ . The statement of the Lemma then follows from the fact that  $\|f_0 - \widehat{f}^{\mathcal{S}}\|_n^2 > \inf_{\mathcal{S}_1 \subset \mathcal{S}_0} \delta_n^{\mathcal{S}_1} - 2 \sup_{\mathcal{S} \not\supseteq \mathcal{S}_0} \rho_n^{\mathcal{S}} > \inf_{i \in \mathcal{S}_0} \delta_n^{\mathcal{S}_0 \setminus i} - M\varepsilon > M\varepsilon$ .

### S.1.2. Proof of Theorem 4.2

We introduce some more notation. We denote with  $\mathcal{F}_{\mathcal{S}} = \bigcup_{K=1}^n \mathcal{F}_{\mathcal{S}}(K)$  all valid trees that split on directions inside  $\mathcal{S}$  and we write  $\Pi_{K, \mathcal{S}}(\cdot)$  for the conditional prior, given  $K$  and  $\mathcal{S}$ .

Similarly as in Section S.1.1, we verify the three conditions (Prior Concentration, Entropy, Prior Anti-concentration). The prior model concentration condition is again satisfied automatically from the definition of model weights in (20) and  $K_{\mathcal{S}_0} = \lfloor C_K / C_\varepsilon n \varepsilon_{n, \mathcal{S}_0}^2 / \log n \rfloor$ . Namely,

$$\pi(K_{\mathcal{S}_0}, \mathcal{S}_0) \propto e^{-C \max\{C_K / C_\varepsilon n \varepsilon_{n, \mathcal{S}_0}^2, q_0 \log p\}} \geq e^{-n \varepsilon_{n, \mathcal{S}_0}^2}, \quad (34)$$

for  $C_K < C_\varepsilon / C$ , where we used the assumption  $q_0 \log p \leq n^{q_0 / (2\alpha + q_0)}$ . Next, the prior concentration in the parameter space associated with the true model

$$\Pi_{K_{\mathcal{S}_0}, \mathcal{S}_0}(f_{\mathcal{T}, \beta} \in \mathcal{F}_{\mathcal{S}_0}(K_{\mathcal{S}_0}) : \|f_{\mathcal{T}, \beta} - f_0\|_n \leq \varepsilon_{n, \mathcal{S}_0}) \geq e^{-dn \varepsilon_{n, \mathcal{S}_0}^2}$$

follows again from Section 8.2 of RP17.

For the entropy considerations, we focus only on models with up to  $q_n$  covariates and up to  $K_n$  splits, where  $q_n = \lceil C_q n \varepsilon_{n, \mathcal{S}_0}^2 / \log p \rceil$  and  $K_n = \lceil \bar{C} n \varepsilon_{n, \mathcal{S}_0}^2 / \log n \rceil$  were defined in Theorem 4.2. This restriction is justified by the following Lemma.

LEMMA S.1.3. *Denote with  $q_n = \lceil C_q n \varepsilon_{n, \mathcal{S}_0}^2 / \log p \rceil$  and  $K_n = \lceil \bar{C} n \varepsilon_{n, \mathcal{S}_0}^2 / \log p \rceil$ . Under the assumptions of Theorem 4.1, we have*

$$\Pi(q \geq q_n | \mathbf{Y}^{(n)}) \rightarrow 0 \quad \text{and} \quad \Pi(K \geq K_n | \mathbf{Y}^{(n)}) \rightarrow 0 \quad (35)$$

in  $\mathbb{P}_{f_0}^{(n)}$ -probability as  $n \rightarrow \infty$ .

PROOF. It suffices to show that  $\Pi(q > q_n) e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \rightarrow 0$  and  $\Pi(K \geq K_n) e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \rightarrow 0$  for  $d > 2$  from (27). We have  $q_0 \leq q_n$  for  $n$  large enough, since  $q_0 = \mathcal{O}(1)$  as  $n \rightarrow \infty$ , and thereby

$$\begin{aligned} \Pi(q \geq q_n) e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} &\lesssim e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \sum_{q=q_n}^p \binom{p}{q} \sum_{K=1}^n e^{-C \max\{K \log n, q \log p\}} \\ &\leq e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \sum_{q=q_n}^p e^{\log n + q \log(p e/q) - C q \log p} \leq e^{\log p + \log n - (C-1) q_n \log p + (d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \\ &\leq e^{-(C-3) q_n \log p + (d+2)n \varepsilon_{n, \mathcal{S}_0}^2}, \end{aligned}$$

where we used the fact that for  $q_0 \geq 2$  and  $\alpha \in (0, 1]$ , we have  $\log n \leq n^{q_0/(2\alpha+q_0)}$ . Since  $q_n \log p \geq C_q n \varepsilon_{n, \mathcal{S}_0}^2$ , the right hand side above goes to zero when  $(C-3)C_q > d+2$ . This will be guaranteed with  $C > 3$  and  $C_q$  large enough. Similarly, we have

$$\begin{aligned} \Pi(K \geq K_n) e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} &\lesssim e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \sum_{q=0}^p \binom{p}{q} \sum_{K=K_n}^n e^{-C \max\{K \log n, q \log p\}} \\ &\leq e^{(d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \sum_{q=0}^p \sum_{K=K_n}^n e^{-(C-1) \max\{K \log n, q \log p\}} \\ &\leq e^{\log(p+1) + \log n - (C-1) K_n \log n + (d+2)n \varepsilon_{n, \mathcal{S}_0}^2} \leq e^{-(C-2) K_n \log n + (d+3)n \varepsilon_{n, \mathcal{S}_0}^2}, \end{aligned}$$

where we used our assumption  $\log p \leq n^{q_0/(2\alpha+q_0)}$ . Since  $K_n \geq \bar{C} n \varepsilon_{n, \mathcal{S}_0}^2$ , the right hand side above goes to zero when  $(C-2)\bar{C} > d+3$ . Together with the prior mass conditions (27) and (34), (35) follows from Lemma 1 of Ghosal and van der Vaart (2007).

This Lemma essentially says that the posterior does not overfit in terms of both  $q$  and  $K$ , where the mass concentrates on models with  $K < K_n$  splits. Note that  $K_n$  is of the same order as the optimal regularity  $K_{\mathcal{S}_0}$ . Now, we denote with  $\Gamma_n \subset \Gamma$  a sieve consisting of all models with less than  $q_n$  variables and  $K_n$  splits. For the entropy bounds of overfitting and underfitting models (inside the sieve  $\Gamma_n$ ), we can use the same arguments as in Section S.1.1. Assume a model  $(K, \mathcal{S}) \in \Gamma_n$ . Then it follows from (30) that

$$\log N(\varepsilon_{n, \mathcal{S}}; \mathcal{F}_{\mathcal{S}}(K); \|\cdot\|_n) \leq K \log(3 B n^3 n^{3/2}) \lesssim K_n \log n \lesssim n \varepsilon_{n, \mathcal{S}_0}^2.$$

For over-fitting models, this can be further upper-bounded with a multiple of  $n \varepsilon_{n, \mathcal{S}}^2$ , thus satisfying (29). The last requirement for the entropy condition is verifying the following variant of (31)

$$\sum_{(K, \mathcal{S}) \in \Gamma_n: \mathcal{S} \not\supset \mathcal{S}_0 \cup K < K_{\mathcal{S}_0}} e^{-C_2 M^2 n \varepsilon_{n, \mathcal{S}_0}^2} + \sum_{(K, \mathcal{S}) \in \Gamma_n: \mathcal{S} \supset \mathcal{S}_0 \cap K \geq K_{\mathcal{S}_0}} e^{-C_2 n \varepsilon_{n, \mathcal{S}}^2} \leq 1 \quad (36)$$

for some suitable  $C_2 > 0$ . Since  $n \varepsilon_{n, \mathcal{S}_0}^2 \leq n \varepsilon_{n, \mathcal{S}}^2$  for  $\mathcal{S} \supset \mathcal{S}_0$ , we can upper-bound the left hand side with

$$\sum_{\mathcal{S}: |\mathcal{S}| < q_n} \sum_{K=1}^{K_n} e^{-C_2 n \varepsilon_{n, \mathcal{S}_0}^2} \leq e^{-C_2 n \varepsilon_{n, \mathcal{S}_0}^2} \left( \frac{2ep}{q_n} \right)^{q_n+1} e^{\log K_n} \leq e^{-C_2 n \varepsilon_{n, \mathcal{S}_0}^2 + (q_n+1) \log p + \log K_n}. \quad (37)$$

Since  $q_n \log p \asymp n \varepsilon_{n, \mathcal{S}_0}^2$  and  $\log K_n \lesssim n^{q_0/(2\alpha+q_0)} \lesssim n \varepsilon_{n, \mathcal{S}_0}^2$ , the right-hand side of (37) converges to zero for some suitably large  $C_2$  as  $n \rightarrow \infty$ , thus satisfying (36).

In place of the anti-concentration condition (similarly as in (33)), we need to verify that the prior probability of larger models (that contain the truth) is small in the sense that, for some  $H > 0$ ,

$$\sum_{(K, \mathcal{S}) \in \Gamma_n: \{\mathcal{S} \supset \mathcal{S}_0\} \cap \{K \geq K_{\mathcal{S}_0}\}} \pi(\mathcal{S}, K) \leq e^{-H n \varepsilon_{n, \mathcal{S}_0}^2}. \quad (38)$$

We can write

$$\sum_{\mathcal{S} \supset \mathcal{S}_0: |\mathcal{S}| < q_n} \sum_{K=K_{\mathcal{S}_0}}^{K_n} \pi(\mathcal{S}, K) \leq \sum_{q=0}^{q_n} \binom{p}{q} \sum_{K=K_{\mathcal{S}_0}}^{K_n} e^{-C K_{\mathcal{S}_0} \log n} \quad (39)$$

$$\leq \left( \frac{2ep}{q_n} \right)^{q_n+1} e^{\log K_n} e^{-C K_{\mathcal{S}_0} \log n}. \quad (40)$$

Because  $q_n \log p \asymp n \varepsilon_{n, \mathcal{S}_0}^2$  and  $\log K_n \lesssim n^{q_0/(2\alpha+q_0)} \lesssim n \varepsilon_{n, \mathcal{S}_0}^2$  the condition (38) is satisfied for some  $H > 0$  when  $C$  and  $C_K$  are large enough.

### S.1.3. Proof of Theorem 4.3

We modify the notation a bit. We adopt the definition of  $\delta$ -valid ensembles from RP17 (Definition 5.3). With  $\mathcal{F}_{\mathcal{S}}(\mathbf{K})$  we denote all  $\delta$ -valid tree ensembles  $f_{\mathcal{E}, \mathbf{B}}$  that (a) are uniformly bounded (i.e.  $\|f_{\mathcal{E}, \mathbf{B}}\|_{\infty} \leq B$  for some  $B > 0$ ), (b) consist of  $T$  trees with  $\mathbf{K} = (K^1, \dots, K^T)' \in \mathbb{N}^T$  leaves and (c) that split along directions  $\mathcal{S}$ .

We start by showing that the prior model concentration condition is satisfied. From our assumption  $q_0 \log p \leq n^{q_0/(2\alpha+q_0)}$  and definition  $K_{\mathcal{S}_0} < C_K/C_{\varepsilon}^2 n \varepsilon_{n, \mathcal{S}_0}^2 / \log n$  and using (23) and (24), we obtain

$$\pi(\mathcal{S}_0, E(K_{\mathcal{S}_0})) \propto \sum_{T=1}^{K_{\mathcal{S}_0}} e^{-C_T T} \sum_{\mathbf{K} \in \mathbb{N}^T: \sum_{t=1}^T K^t = K_{\mathcal{S}_0}} e^{-C n^{q_0/(2\alpha+q_0)} \log n} \geq e^{-(C_T C_K / (C_{\varepsilon} \log n) + C / C_{\varepsilon}^2) n \varepsilon_{n, \mathcal{S}_0}^2}.$$

The right-hand side can be further lower-bounded with  $e^{-n \varepsilon_{n, \mathcal{S}_0}^2}$  for a large enough  $C_{\varepsilon}$  and  $n$ . Next, we need to show prior concentration in the parameter space under the true model equivalence class  $(\mathcal{S}_0, E(K_{\mathcal{S}_0}))$ . All that is needed is finding a single well-approximating forest supported on one partition ensemble characterized by  $(T, \mathbf{K})$  from the equivalence class  $E(K_{\mathcal{S}_0})$ . Such an ensemble can be obtained by considering  $T = 1$  and a single  $k$ - $d$  tree with  $K_{\mathcal{S}_0}$  leaves from Lemma 3.2 of RP17. The prior concentration condition then boils down to (27), which has already been verified in RP17.

Next, we show that for  $K_n = \lceil C n \varepsilon_{n, \mathcal{S}_0}^2 / \log n \rceil$  we have

$$\Pi \left( (T, \mathbf{K}) : \sum_{t=1}^T K^t \geq K_n \mid \mathbf{Y}^{(n)} \right) \rightarrow 0.$$

We can write

$$\begin{aligned} \Pi \left( (T, \mathbf{K}) : \sum_{t=1}^T K^t \geq K_n \right) &\lesssim \sum_{T=1}^n e^{-C_T T} \sum_{q=0}^p \binom{p}{q} \sum_{Z=K_n}^n \sum_{\mathbf{K}: \sum_{t=1}^T K^t = Z} e^{-C \max\{Z \log n, q \log p\}} \\ &\lesssim e^{-(C-1)K_n \log n + \log p + 2 \log n + \log p(n) - C_T}, \end{aligned}$$

where  $p(n)$  is the partitioning number. According to Andrews (1976), we have

$$\log p(n) \sim \pi \sqrt{\frac{2n}{3}} \quad \text{as } n \rightarrow \infty. \quad (41)$$

Under our assumptions  $q_0 > 2$  and  $\alpha \in (0, 1]$ , we have  $\sqrt{n} \leq n^{q_0/(2\alpha+q_0)}$  and  $\log n \leq n^{q_0/(2\alpha+q_0)}$ . From  $\log p \leq n^{q_0/(2\alpha+q_0)}$  and using the fact that  $K_n \geq \bar{C} n \varepsilon_{n, S_0}^2 / \log n$ , we can then write

$$\Pi \left( (T, \mathbf{K}) : \sum_{t=1}^T K^t \geq K_n \right) e^{(d+2)n \varepsilon_{n, S_0}^2} \lesssim e^{-[(C-1)\bar{C} - D\pi\sqrt{2/3} - d - 5]n \varepsilon_{n, S_0}^2}$$

for some  $D > 0$ . The right hand side goes to zero for  $C > 1$  and  $\bar{C}$  large enough. Similarly, we can show that  $\Pi(q \geq q_n | \mathbf{Y}^{(n)}) \rightarrow 0$  as  $n \rightarrow \infty$  for  $q_n = \lceil C_q n \varepsilon_{n, S_0}^2 / \log p \rceil$  by proceeding as in Lemma S.1.3 in Section S.1.2.

Based on the previous paragraph, we narrow down attention to a subset of model indices  $\Gamma_n \subset \Gamma$ , consisting of models  $(\mathcal{S}, E(Z))$  such that  $|\mathcal{S}| < q_n$  and  $Z < K_n$ . We now define a sieve  $\mathcal{F}_n$  as follows

$$\mathcal{F}_n = \bigcup_{q=0}^{q_n} \bigcup_{T=1}^{K_n} \bigcup_{\sum_{t=1}^T K^t \leq K_n} \bigcup_{\mathcal{S}: |\mathcal{S}|=q} \mathcal{F}_{\mathcal{S}}(\mathbf{K}).$$

It follows from the previous paragraph that  $\Pi(\mathcal{F}_n^c | \mathbf{Y}^{(n)}) \rightarrow 0$  as  $n \rightarrow \infty$ . For the entropy calculation we thus focus on the sieve  $\mathcal{F}_n$ .

We first note that the metric entropy  $\log N(\varepsilon_{n, \mathcal{S}}; \mathcal{F}(\mathcal{E}); \|\cdot\|_n)$ , where  $\mathcal{F}(\mathcal{E})$  are all uniformly bounded forests supported on a  $\delta$ -valid partition ensemble  $\mathcal{E}$ , can be upper-bounded with  $\left(\sum_{t=1}^T K^t\right) \log(B/\varepsilon_{n, \mathcal{S}} C_1 \kappa(\mathcal{E}) \sqrt{n})$  (follows from equation (9.3) of RP17), where  $\kappa(\mathcal{E})$  is the condition number of a valid ensemble (defined in Section 9.1. of RP17). Next, we find an upper bound for the covering number of the tree ensembles that are attached to a model  $(\mathcal{S}, E(Z))$ , where  $E(Z)$  is the equivalence class of  $(T, \mathbf{K})$  defined in (22). From Section 9.1 of RP17, and using the fact that  $\Delta(E(Z)) \leq Z!p(Z)$ , it follows that

$$\begin{aligned} & \log N \left( \varepsilon_{n, \mathcal{S}}; \bigcup_{(T, \mathbf{K}) \in E(Z)} \mathcal{F}_{\mathcal{S}}(\mathbf{K}) \cap \mathcal{F}_n; \|\cdot\|_n \right) \\ & \leq \log \Delta(E(Z)) + \log \Delta(\mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}}) + Z \log(B/\varepsilon_{n, \mathcal{S}} C_1 \kappa(\mathcal{E}) \sqrt{n}) \\ & \lesssim Z \log Z + \sqrt{Z} + Z \log(|\mathcal{S}|n^2) + Z \log \left( n^{2+\delta/2} \sqrt{Z} \right) \end{aligned}$$

for some  $C_1 > 0$ , where  $\Delta(\mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}})$  is the cardinality of  $\delta$ -valid ensembles  $\mathcal{V}\mathcal{E}_{\mathcal{S}}^{\mathbf{K}}$ . Inside the sieve, we have  $|\mathcal{S}| < q_n \leq n$  and  $Z < K_n \asymp n \varepsilon_{n, S_0}^2 / \log n$  and thereby we can upper bound the log entropy with a constant multiple of  $n \varepsilon_{n, S_0}^2$ . For an overfitting model  $(\mathcal{S}, E(Z))$  such that  $Z \geq K(\mathcal{S}_0)$  and  $\mathcal{S} \supset \mathcal{S}_0$ , the log-covering number is further upper-bounded with  $n \varepsilon_{n, \mathcal{S}}^2 \geq n \varepsilon_{n, S_0}^2$ . Next, we verify the following variant of condition (31)

$$\sum_{\Gamma_n \cap \Gamma_{\{\mathcal{S} \not\supset \mathcal{S}_0\} \cup \{Z < K_{\mathcal{S}_0}\}}} e^{-C_2 M^2 n \varepsilon_{n, S_0}^2} + \sum_{\Gamma_n \cap \Gamma_{\{\mathcal{S} \supset \mathcal{S}_0\} \cap \{Z \geq K_{\mathcal{S}_0}\}}} e^{-C_2 n \varepsilon_{n, \mathcal{S}}^2} \leq 1 \quad (42)$$

for some  $C_2 > 0$ . Since  $n\varepsilon_{n,\mathcal{S}}^2 > n\varepsilon_{n,\mathcal{S}_0}^2$  for  $\mathcal{S} \supset \mathcal{S}_0$  and  $M > 1$ , we can upper-bound the left-hand-side with

$$e^{-C_2 n\varepsilon_{n,\mathcal{S}_0}^2} \sum_{q=0}^{q_n} \binom{p}{q} \sum_{Z=1}^{K_n} \Delta(E(Z)) \lesssim \left(\frac{2ep}{q_n}\right)^{q_n+1} e^{-C_2 n\varepsilon_{n,\mathcal{S}_0}^2 + \log q_n + \log K_n + K_n \log K_n + \pi\sqrt{2K_n/3}},$$

where we used the fact  $\Delta(E(Z)) \leq Z!p(Z)$  and (41). Since  $K_n \log K_n \lesssim n\varepsilon_{n,\mathcal{S}_0}^2$  and  $q_n \log p \asymp n\varepsilon_{n,\mathcal{S}_0}^2$ , the right hand side goes to zero for a large enough constant  $C_2 > 0$ .

Lastly, the anti-concentration condition is replaced with

$$\sum_{T=K_n}^n \pi(T) \sum_{\Gamma_n \cap \Gamma_{\{\mathcal{S} \supset \mathcal{S}_0\} \cap \{Z \geq K_{\mathcal{S}_0}\}}} \sum_{\mathbf{K} \in \mathbb{N}^T: \sum_{t=1}^T K^t = Z} \pi(\mathcal{S}, \mathbf{K} | T) \leq e^{-H n\varepsilon_{n,\mathcal{S}_0}^2}$$

for some  $H > 0$ . Using the fact  $\pi(\mathcal{S}, \mathbf{K} | T) \gtrsim e^{-C \sum K^t \log n}$ , we can upper-bound the left hand side above with

$$\begin{aligned} & \sum_{T=1}^{K_n} \pi(T) e^{-C K_{\mathcal{S}_0} \log n} \sum_{\Gamma_n \cap \Gamma_{\{\mathcal{S} \supset \mathcal{S}_0\} \cap \{Z \geq K_{\mathcal{S}_0}\}}} \Delta(E(Z)) \\ & \lesssim e^{-C K_{\mathcal{S}_0} \log n} \left(\frac{2ep}{q_n}\right)^{q_n+1} e^{2 \log K_n + K_n \log K_n + \pi\sqrt{2K_n/3} - C_T} \end{aligned}$$

Using similar arguments as before, and because  $K_{\mathcal{S}_0} \log n \geq C_K/C_\varepsilon n\varepsilon_{n,\mathcal{S}_0}^2$ , the condition will be satisfied for large enough  $C > 0$  and  $C_K > 0$ .

#### S.1.4. Theory for ABC

First, we show the following ABC posterior concentration result.

**THEOREM S.1.1.** *Under the assumptions of Theorem 4.1 and assuming  $\sigma^2 = 1/n$  in (1), the naive ABC posterior satisfies with  $\mathbb{P}_{f_0}^{(n)}$  tending to one*

$$\Pi [\|f - f_0\|_n > \lambda_n \mid \|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] \lesssim 1/M$$

for  $\epsilon_n^T = \sqrt{2 \log n/n}$ ,  $\lambda_n = 4\epsilon_n^T/3 + 1/\sqrt{n}$  and for any  $M > 0$  large enough.

**PROOF.** We will be working conditionally on the event  $\mathcal{A} = \{\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)' : \max_{1 \leq i \leq n} |\varepsilon_i| \leq \sqrt{2 \log n/n}\}$  whose complement has a small probability, i.e.  $\mathbb{P}_{f_0}^{(n)}[\mathcal{A}^c] \leq c_0/\sqrt{2 \log n}$  for some  $c_0 > 0$  when  $\varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, 1/n)$ . On the event  $\mathcal{A}$ , we have

$$\|\mathbf{Y} - f_0\|_n = \sqrt{\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2} \leq \sqrt{2 \log n/n} \equiv \epsilon_n^T.$$

We now define a joint event

$$\mathcal{A}(\epsilon_n^T, \lambda_n) \equiv \{(\mathbf{Y}^*, f) : \|\mathbf{Y}^* - \mathbf{Y}\|_n \leq \epsilon_n^T \quad \text{and} \quad \|f - f_0\|_n > \lambda_n\}.$$



For all  $(\mathbf{Y}^*, f) \in \mathcal{A}(\epsilon_n^T, \lambda_n)$  we have

$$\|f - f_0\|_n \leq \|\mathbf{Y}^* - \mathbf{Y}\|_n + \|f - \mathbf{Y}^*\|_n + \|f_0 - \mathbf{Y}\|_n \leq \frac{4}{3}\epsilon_n^T + \|f - \mathbf{Y}^*\|_n.$$

This means that  $(\mathbf{Y}^*, f) \in \mathcal{A}(\epsilon_n^T, \lambda_n)$  implies  $\|f - \mathbf{Y}^*\|_n > \lambda_n - \frac{4}{3}\epsilon_n^T$  and choosing  $\lambda_n \geq \frac{4}{3}\epsilon_n^T + t_\varepsilon$  leads to

$$\Pi[\mathcal{A}(\epsilon_n^T, \lambda_n)] \leq \int \mathbb{P}_f[\|f - \mathbf{Y}^*\|_n > t_\varepsilon] d\Pi(f)$$

and

$$\Pi \left[ \|f - f_0\|_n > \frac{4}{3}\epsilon_n^T + t_\varepsilon \mid \|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T \right] \leq \frac{\int \mathbb{P}_f[\|\mathbf{Y}^* - f\|_n > t_\varepsilon] d\Pi(f)}{\int \mathbb{P}_f[\|\mathbf{Y}^* - \mathbf{Y}\|_n \leq \epsilon_n^T] d\Pi(f)}. \quad (43)$$

Now, we have for a random variable  $\chi_n^2$  with a chi-square distribution with  $n$  degrees of freedom

$$\mathbb{P}_f[\|\mathbf{Y}^* - f\|_n > u] = \mathbb{P}_f \left[ \frac{\chi_n^2}{n^2} > u^2 \right] = \mathbb{P}_f \left[ e^{\chi_n^2/4} > e^{u^2 n^2/4} \right] \leq \frac{2^{n/2}}{e^{u^2 n^2/4}}.$$

Next, for  $n$  large enough we can write

$$\int \mathbb{P}_f[\|\mathbf{Y}^* - \mathbf{Y}\|_n \leq \epsilon_n^T] d\Pi(f) \geq \int_{\|f - f_0\|_n \leq \epsilon_n^T/3} \mathbb{P}_f[\|\mathbf{Y}^* - f\|_n \leq \epsilon_n^T/3] d\Pi(f) \quad (44)$$

$$\geq \Pi[\|f - f_0\|_n \leq \epsilon_n^T/3] - e^{n/2 \log 2 - n \log n/18} \quad (45)$$

$$\geq \Pi[\|f - f_0\|_n \leq \epsilon_n^T/3]/2. \quad (46)$$

Next (under the assumption  $q_0 \log p < n^{q_0/(2\alpha+q_0)}$ , we have  $\pi(\mathcal{S}_0) \geq e^{-n\varepsilon_{n,s_0}^2}$  and (assuming  $K = K_{\mathcal{S}_0} \asymp n\varepsilon_{n,s_0}^2/\log n$  and denoting  $\widehat{\beta} \in \mathbb{R}^K$  the steps of the  $\|\cdot\|_n$  projection of  $f_0$  onto trees with  $K$  leafs) for some  $c > 0$

$$\Pi[\|f - f_0\|_n \leq \epsilon_n^T/3] > e^{-n\varepsilon_{n,s_0}^2} \Pi(\|\beta - \widehat{\beta}\|_2 \leq \epsilon_n^T/6) \quad (47)$$

$$> e^{-n\varepsilon_{n,s_0}^2} \frac{e^{-K \log 2 - \|\widehat{\beta}\|_2^2 - (\epsilon_n^T)^2/72 + K/2 \log[(\epsilon_n^T)^2/36]}}{\Gamma(K/2)K/2} > e^{-cn\varepsilon_{n,s_0}^2}. \quad (48)$$

We can now upper-bound (43) with  $2^{n/2} e^{-t_\varepsilon^2 n^2/4 + cn\varepsilon_{n,s_0}^2}$  which is smaller than an arbitrary constant  $M > 0$  for  $n$  large enough if we choose  $t_\varepsilon = 1/\sqrt{n}$ .

Given this consistency result, we can immediately conclude (using the inequality in (21) in the paper) that the ABC posterior will not reward underfitting model as long as our identifiability and irrepresentability conditions are satisfied with  $\varepsilon = \lambda_n$ . In other words, under the assumptions of Theorem S.1.1 and assuming that  $\mathcal{S}_0$  is  $(f_0, \lambda_n)$ -identifiable and that  $\lambda_n$ -irrepresentability holds we have, with  $\mathbb{P}_{f_0}$  tending to one and for any  $M > 0$ ,

$$\Pi[\mathcal{S} \not\supseteq \mathcal{S}_0 \mid \|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] \lesssim 1/M.$$

Regarding over-fitting models, we first show the following ABC analogue of Lemma 8.1. We can write, on the event  $\mathcal{A}$ , and for  $q_n = C_q [n \epsilon_{n, \mathcal{S}_0}^2 / \log p]$  (as in Lemma 8.1)

$$\Pi_1 \equiv \Pi [q \geq q_n \mid \|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] = \sum_{\mathcal{S}: |\mathcal{S}| \geq q_n} \pi(\mathcal{S}) \frac{\int \mathbb{P}_f[\|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] d\Pi(f \mid \mathcal{S})}{\int \mathbb{P}_f[\|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] d\Pi(f)}.$$

It turns out from the proof of Theorem S.1.1 that

$$\Pi_1 \leq \frac{\sum_{q \geq q_n} \sum_{\mathcal{S}: |\mathcal{S}|=q} \pi(\mathcal{S})}{\int \mathbb{P}_f[\|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] d\Pi(f)} \lesssim e^{c n \epsilon_{n, \mathcal{S}_0}^2} \Pi(q \geq q_n).$$

In the proof of Lemma 1.1 we have already showed (under the assumptions of Theorem 4.1) that  $\Pi(q \geq q_n) \lesssim e^{-n \epsilon_{n, \mathcal{S}_0}^2 C}$  for some  $C > 0$ . Choosing  $C_q$  large enough, one concludes that  $\Pi_1 \rightarrow 0$  as  $n \rightarrow \infty$ . This shows that the ABC posterior concentrates on the sieve of models  $\mathcal{F}_n$  with up to  $q_n$  covariates. Using this result, we can focus on models of size up to  $q_n$  and show that the posterior probability of over-fitting models goes to zero. Indeed, on the event  $\mathcal{A}$  and on  $\mathcal{F}_n$  we have (using an inequalities (4) and (6))

$$\begin{aligned} \Pi [\{\mathcal{S} \supset \mathcal{S}_0\} \cap \mathcal{F}_n \mid \|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] &\leq \frac{\sum_{\mathcal{S} \supset \mathcal{S}_0: |\mathcal{S}| \leq q_n} \pi(\mathcal{S})}{\int \mathbb{P}_f[\|\mathbf{Y} - \mathbf{Y}^*\|_n \leq \epsilon_n^T] d\Pi(f)} \\ &\lesssim e^{(c-C) n \epsilon_{n, \mathcal{S}_0}^2} \left(\frac{2ep}{q_n}\right)^{q_n+1} \lesssim e^{-H n \epsilon_{n, \mathcal{S}_0}^2} \end{aligned}$$

for some  $H > 0$  with  $C > 0$  is large enough. This concludes that the ABC posterior will lead to consistent variable selection as well.

We now discuss how the theory can be extended when data-splitting is deployed in ABC. First, we discuss the case when the split is done only once before applying ABC (not internally at each iteration). Denote with  $n_1$  the training sample size and with  $n_2$  the validation sample size. In order for the consistency result in Theorem S.1.1 to hold, we need to make sure that prior concentration holds in the sense that  $\Pi[\|f - f_0\|_{n_2} \lesssim \epsilon_{n_2}^T] \geq e^{-c n_2 \epsilon_{n_2, \mathcal{S}_0}^2}$  for some  $c > 0$ . Leaving  $n_1$  data-points for training the prior, we know (from results in RP17 under fixed  $\sigma^2$ ) that the posterior concentrates at the optimal rate (up to a log factor), i.e.

$$\Pi[\|f - f_0\|_{n_1} \lesssim \epsilon_{n_1, \mathcal{S}_0} \mid \mathbf{Y}_{\mathcal{I}}^{(n)}, \mathcal{S}_0] \rightarrow 1 \quad \text{as } n_1 \rightarrow \infty.$$

Choosing  $n_1$  and  $n_2$  in such a way so that  $\epsilon_{n_1, \mathcal{S}_0} \lesssim \epsilon_{n_2}^T \equiv \sqrt{2(\log n_2)/n_2}$  (and assuming that observed fixed covariates in the training and testing sets are close), the prior concentration condition will be satisfied and the ABC will be consistent and concentrate at the rate  $\lambda_{n_2}$ . This implies variable selection consistency of our ABC method under identifiability and irrepresentability conditions which depend on  $\lambda_{n_2}$ . A similar conclusion is obtained for the expected posterior prior (9) where

$$\Pi[\|f - f_0\|_{n_1} \lesssim \epsilon_{n_1, \mathcal{S}_0}] \geq \pi(\mathcal{S}_0) \frac{1}{L} \sum_l \Pi[\|f - f_0\|_{n_1} \lesssim \epsilon_{n_1, \mathcal{S}_0} \mid \mathbf{Y}_{\mathcal{I}_l}^{(n)}, \mathcal{S}_0] \gtrsim \pi(\mathcal{S}_0).$$

A rigorous proof of ABC consistency for the expected posterior priors would require more care and will be left for future investigation.

**Table S1.** Computation time of 1 000 MCMC iterations of BART/DART in seconds (using the Friedman’s datasets with  $\sigma = 5$  and autocorrelation of 0.9).

		BART				DART			
		$T = 10$	$T = 20$	$T = 50$	$T = 200$	$T = 10$	$T = 20$	$T = 50$	$T = 200$
$n = 100$	$p = 100$	0.21	0.32	0.54	1.86	0.55	0.64	0.76	2.27
	$p = 1000$	0.53	0.67	1.57	5.48	0.96	0.98	1.91	5.79
	$p = 10000$	3.56	5.58	10.91	39.16	5.93	7.72	12.55	36.95
$n = 250$	$p = 100$	0.21	0.34	0.79	2.99	0.41	0.56	0.99	3.19
	$p = 1000$	0.50	0.82	1.81	6.51	0.87	1.14	1.83	6.58
	$p = 10000$	3.70	5.63	11.38	40.48	6.29	8.06	12.97	39.13
$n = 500$	$p = 100$	0.29	0.53	1.23	4.93	0.49	0.71	1.40	5.07
	$p = 1000$	0.63	1.11	2.36	8.65	1.01	1.30	2.29	7.89
	$p = 10000$	4.21	6.54	12.27	43.35	6.80	8.35	13.46	37.81
$n = 1000$	$p = 100$	0.53	0.97	2.22	8.86	0.71	1.13	2.30	8.95
	$p = 1000$	0.91	1.49	3.32	12.82	1.24	1.84	3.54	12.12
	$p = 10000$	5.41	8.18	14.23	48.92	7.61	9.06	14.47	41.90
$n = 10000$	$p = 100$	7.17	10.78	23.25	82.45	6.37	12.07	22.33	92.23
	$p = 1000$	13.00	22.12	34.67	125.98	12.76	16.95	40.25	102.72
	$p = 10000$	25.35	31.39	59.71	218.08	28.59	39.93	73.99	171.73

## S.2. ABC Computational Feasibility

Regarding computational considerations, our sampling method deploys MCMC inside each ABC iteration but uses only on a subset of the original observations (say  $\frac{n}{2}$  observations) and a subset of  $|\mathcal{S}| < p$  variables. In addition, we only need to collect one posterior sample after a burnin period  $B$ .

In order to understand how ABC scales with  $|\mathcal{S}|$ ,  $p$  and  $s$ , we first assess the computing time of plain BART/DART. The timing comparisons are summarized in Table S1. From these computations we can conclude, for example, that running  $M = 1000$  BART iterations with  $T = 200$  trees (the default) on a dataset with  $p = 10000$  variables and  $n = 500$  observations takes 43.35 seconds which roughly amounts to running  $43.35 \times 5/0.5 = 433.5$  ABC iterations with  $B = 200$  burnin MCMC iterations,  $T = 10$  trees and with  $s = n/2$ , assuming that the sparsity prior is such that  $|\mathcal{S}| \approx 1000$ . Under the same settings but a stricter sparsity prior such that  $|\mathcal{S}| \approx 100$ , we obtain  $43.35 \times 5/0.21 = 1032.14$  ABC iterations for the same time as 1000 BART iterations. These computing times, however, do not take into account autocorrelation in BART samples, where  $M = 1000$  BART MCMC iterations do not necessarily yield 1000 *effective* samples. One advantage of ABC sampling over MCMC is that it is embarrassingly parallel and that it does not incur correlation. This provides an opportunity for large speedups using parallel computing.

## S.3. Spike-and-Forests: MCMC Variant

As a precursor to ABC Bayesian Forests, we first implemented an MCMC algorithm for joint sampling from a posterior  $\Pi(\mathcal{S}, \mathcal{E} | \mathbf{Y}^{(n)})$  over the space of models and tree ensemble partitions. We refer to this algorithm as Spike-and-Forests. The sampling follows a Metropolis-Hasting scheme, exploiting the additive structure of forests by sampling each tree individually from conditionals in a Gibbs manner within each Metropolis step (Bayesian backfitting by Chipman et al. (2010)). The key is assigning a joint proposal

distribution  $pr(\mathcal{S}, \mathcal{E} | \mathcal{S}_m, \mathcal{E}_m) = pr(\mathcal{S} | \mathcal{S}_m)pr(\mathcal{E} | \mathcal{S}, \mathcal{E}_m)$  over variable subsets  $\mathcal{S}$  and partition ensembles  $\mathcal{E}$ , where  $\mathcal{S}_m$  and  $\mathcal{E}_m$  are current MCMC states.

We explain the proposal mechanism using a single tree and write  $\mathcal{T}$  instead of  $\mathcal{E}$ . First, a model proposal  $\mathcal{S}^*$  is sampled from  $pr(\mathcal{S} | \mathcal{S}_m)$  which consists of the following three options: **add**, **delete** and **stay** for adding/deleting one (or none) of the variables. These three steps are chosen with probabilities 0.4, 0.4 and 0.2, respectively. Candidate variables for deletion/addition are chosen from a uniform distribution. Given the newly suggested model  $\mathcal{S}^*$ , the proposal distribution  $pr(\mathcal{T} | \mathcal{S}^*, \mathcal{T}_m)$  consists of various moves, described below, depending on the status of  $\mathcal{S}^*$ .

If  $\mathcal{S}^*$  was obtained from  $\mathcal{S}_m$  by adding a variable, the proposal  $pr(\mathcal{T} | \mathcal{S}^* = \text{add}, \mathcal{T}_m)$  consists of two steps: **birth** and **replace**. In the **birth** step, a bottom node is added to  $\mathcal{T}_m$  and in the **replace** step one of the variables that occurs more than once inside  $\mathcal{T}_m$  is replaced with the new variable. The birth step increases the size of the tree, while the replace step does not. The two steps are chosen with probabilities

$$\pi_{\text{birth,add}} = 0.7 \min \left\{ \frac{\pi(K+1)}{\pi(K)}, 1 \right\}, \pi_{\text{birth,replace}} = 1 - \pi_{\text{birth,add}},$$

where  $K$  is the number of bottom nodes in  $\mathcal{T}_m$  and  $\pi(K)$  is a prior on the number of bottom nodes. If no variable appears more than once in the tree, then **replace** is invalid and  $\pi_{\text{birth,replace}}$  is set to 0.

If  $\mathcal{S}^*$  is obtained from  $\mathcal{S}_m$  by deleting a variable, the proposal  $pr(\mathcal{T} | \mathcal{S}^* = \text{delete}, \mathcal{T}_m)$  consists of two steps: **death** and **replace**. If the variable chosen for deletion occurs in a bottom node, it can be removed from a tree  $\mathcal{T}_m$  with a **delete** step that erases the bottom node. If the variable occurs inside the tree, it can be deleted by replacing it with other variables in the **replace** step. If both of these moves are eligible, we pick one of them with probabilities

$$\pi_{\text{death,delete}} = 0.7 \min \left\{ \frac{\pi(K-1)}{\pi(K)}, 1 \right\}, \pi_{\text{death,replace}} = 1 - \pi_{\text{death,delete}}.$$

If the variable suggested for deletion is not in a bottom node, then  $\pi_{\text{death,delete}} = 0$ .

If the pool of variables stays the same, i.e.  $\mathcal{S}^* = \mathcal{S}_m$ , the proposal  $pr(\mathcal{T} | \mathcal{S}^* = \text{stay}, \mathcal{T}_m)$  consists of 4 moves: **add**, **delete**, **replace** and **rule**. All proposal moves, and their probabilities, are adopted from Bayesian CART of [Denison et al. \(1998\)](#). These steps only modify the tree configuration without adding/deleting variables.

Regarding the prior distributions for our MCMC implementation, we assume the beta-binomial prior on the variable subsets. Namely, for binary indicators  $\gamma_j \in \{0, 1\}$ , for whether or not  $x_j$  is active, we assume  $\mathbb{P}(\gamma_j = 1 | \theta) = \theta$  and  $\theta \sim \mathcal{B}(a, b)$ . The prior distribution on trees consists of (a) the truncated Poisson distribution on the number of bottom leaves, (b) uniform prior over trees with the same number of leaves and (c) standard Gaussian prior on the step sizes. This is the Bayesian CART prior proposed by [Denison et al. \(1998\)](#) and analyzed theoretically by [Ročková and van der Pas \(2017\)](#). In the computation of MH acceptance ratios, we leverage the fact that the bottom leave parameters can be integrated out to obtain a conditional marginal likelihood, given each partition.

The MCMC sampling routine can be extended to spike-and-forests, altering each tree inside the forests one by one through Bayesian backfitting (Chipman et al., 2010). One big advantage of the Bayesian forest representation is that it accelerates mixing since most trees are shallow and thereby more easily modified throughout MCMC (see Pratola (2016)).

#### S.4. Sensitivity Analysis

Our sensitivity analysis focuses on two aspects. First, we want to assess how the choices of  $M$  (the number of ABC samples),  $T$  (the number of trees in each forest),  $B$  (the number of burn-in iterations inside each ABC iteration) and  $\epsilon$  (tolerance for ABC acceptance) collaboratively impact ABC variable importance. Second, we want to investigate the impact of different data splitting strategies, including varying choices of  $s$  (proportion of data used in training) and pre-determined data splitting versus internal data splitting. There is an obvious tradeoff between  $s$  and  $M$ , where small  $s$  will yield fewer ABC pseudo-observations that are compatible with the observed data and  $M$  will thereby have to be larger. We have considered the following combinations

$$M \in \{1\,000, 10\,000\} \times T \in \{10, 25, 50\} \times B \in \{200, 1\,000\} \times \epsilon \in \{\text{top } 1\%, 5\%, 10\%\}$$

These comparisons are conducted using the Friedman’s simulation setup with  $p \in \{100, 1\,000\}$ ,  $\rho = 0.9$  (autoregressive) and  $\sigma = 5$ , assuming  $s = n/2$  and internal splitting for ABC. We also include various sample sizes  $n \in \{100, 500, 1\,000\}$  for each  $p$ . For each setting, we show ABC inclusion probabilities (ip) for the first 30 variables of which only the first 5 are active (Figure S1). We denote the parameters for each ABC setup by  $T \star B$  where, for example,  $20 \star 200$  means each forest consists of  $T = 20$  trees and uses  $B = 200$  MCMC iterations as a burnin.

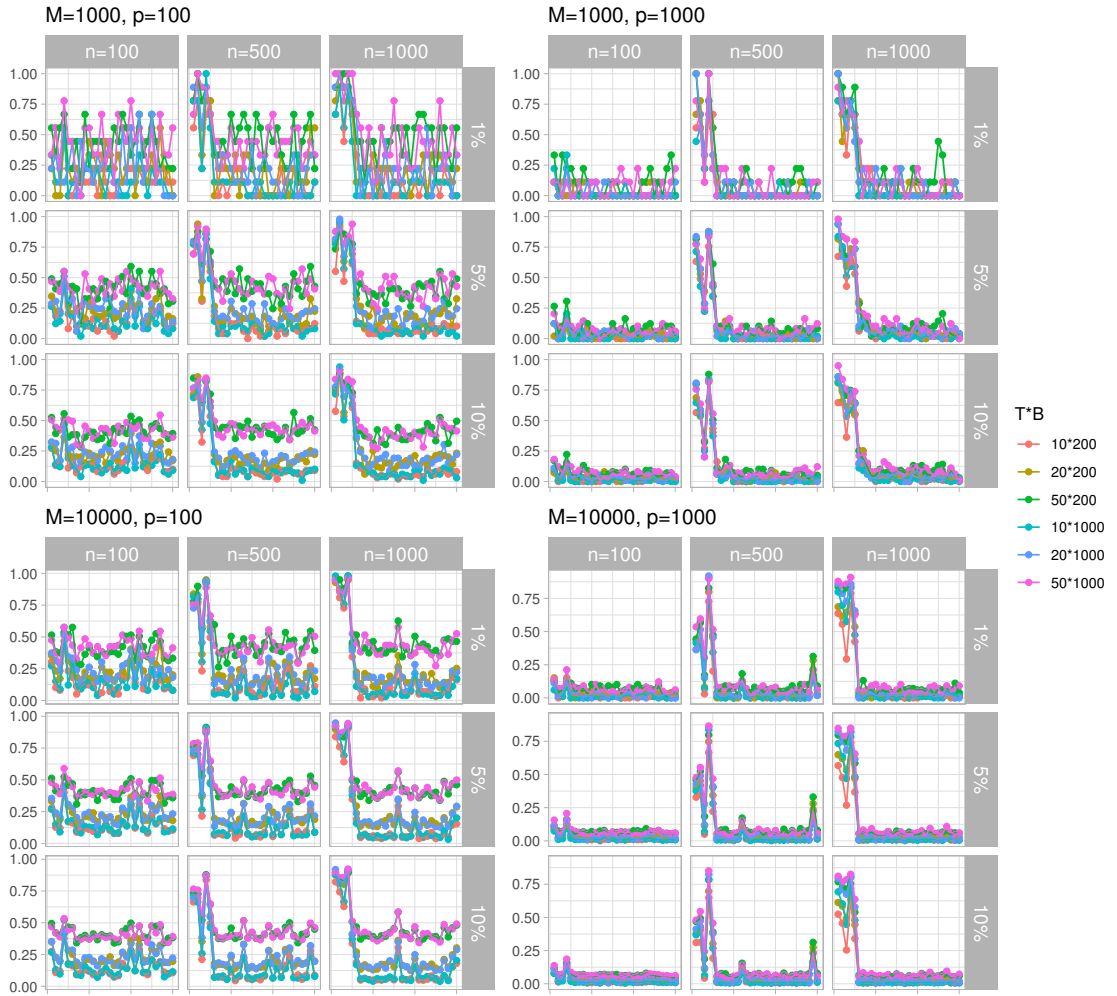
From the figures we can see that ABC is more sensitive to the choice of  $T$  than to the choice of  $B$ . This is not entirely unexpected. As suggested in Chipman et al. (2010) and Bleich et al. (2014), a large value of  $T$  allows for increased flexibility in fitting the model while smaller  $T$  should be adopted for the purpose of variable selection. The variables must compete with each other to be included when  $T$  is small. In terms of a median probability model, the model tends to have more power and higher false discoveries when  $T$  is large, and less power and fewer false discoveries when  $T$  is small.

Regarding  $\epsilon$ , although the trends are similar for top 1%, 5% and 10% selected model, higher variance is observed for smaller tolerance when  $M$  is not large enough, especially for  $M = 1\,000$  with top 1% models accepted. This is, again, not entirely unexpected.

The comparisons in Figure S1 were done assuming  $s = n/2$ . We now consider a similar simulation study, but for  $\epsilon = \{\text{top } 10\%\}$  and various  $s$  by considering

$$M \in \{1\,000, 10\,000\} \times T \in \{10, 25, 50\} \times B \in \{200, 1\,000\} \times s \in \{n/5, n/2, 4n/5\}.$$

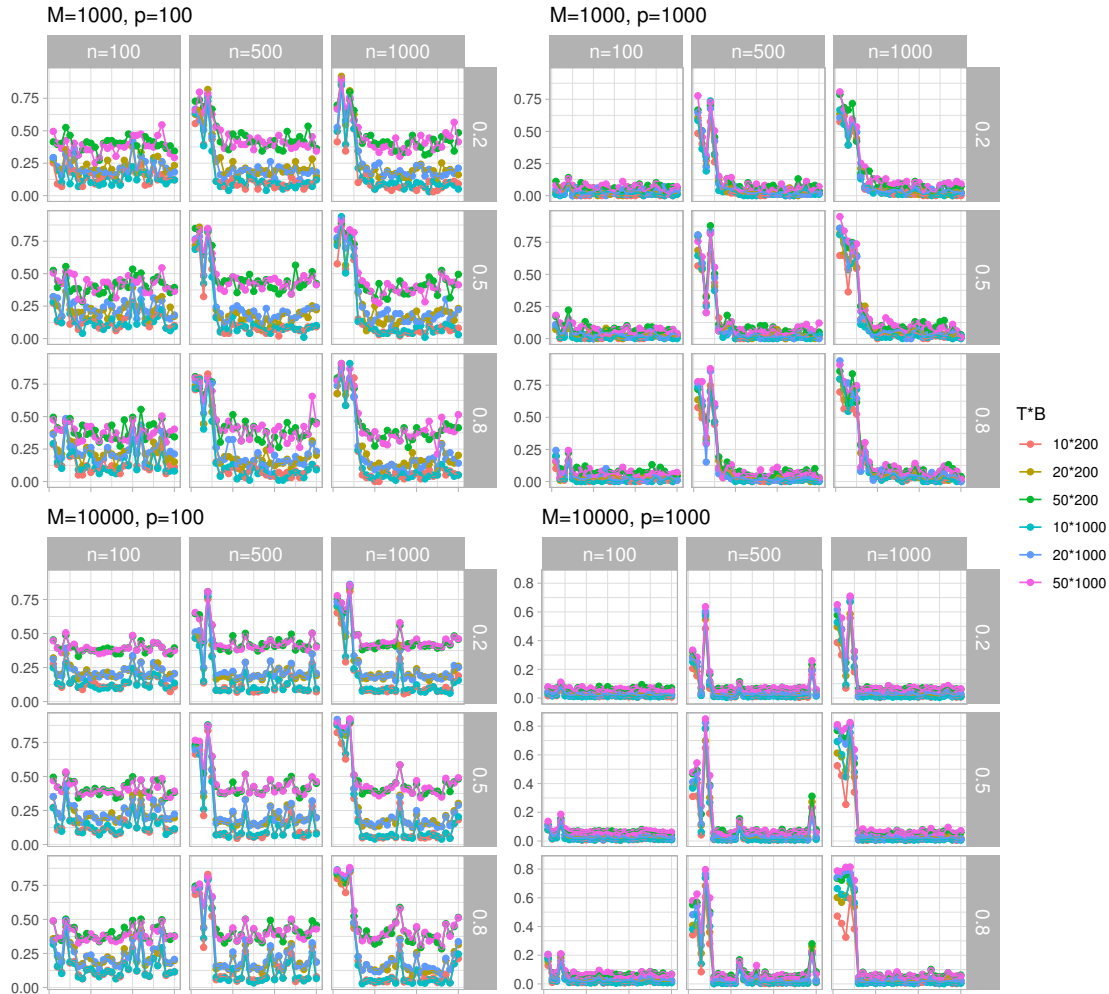
The results are displayed in Figure S2. The posterior inclusion probabilities do not seem to vary much with respect to  $s$ . This suggests that even  $s = 0.2n$  provides reasonable prior guesses for ABC regarding variable selection. Based on this sensitivity analysis, we choose  $T = 20$ ,  $B = 200$ ,  $M = 1\,000$ ,  $s = n/2$ ,  $\epsilon = \{\text{top } 10\%\}$  as the default parameters for our ABC model.



**Fig. S1.** ABC inclusion probabilities of the first 30 variables over different  $\epsilon$ . Each panel corresponds to a different combination of  $p \in \{100, 1000\}$  and  $M \in \{1000, 10000\}$ . Each row indicates a different model averaging strategy based on a different  $\epsilon$  value. Each column corresponds to a different sample size. The legend represents various combinations of  $T \times B$ . For example,  $20 \times 200$  means each forest consists of  $T = 20$  trees and  $B = 200$  MCMC iterations as burnin. Note that we use  $s = n/2$  here.

The last part of the sensitivity analysis we want to investigate the differences between pre-determined data splitting and internal data splitting. Customarily (Berger and Pericchi, 2004), the subsample size  $s$  is chosen as the minimal number of samples needed to convert an improper prior into a proper one. Our situation, however, is different in at least three aspects: (a) we are converting a proper uninformative prior into an informative one, (b) our model is entirely non-parametric and (c) we aim to enhance ABC acceptance rate rather than using non-informative priors for model selection with Bayes factors. As pointed out in Berger and Pericchi (2004), defining any opti-



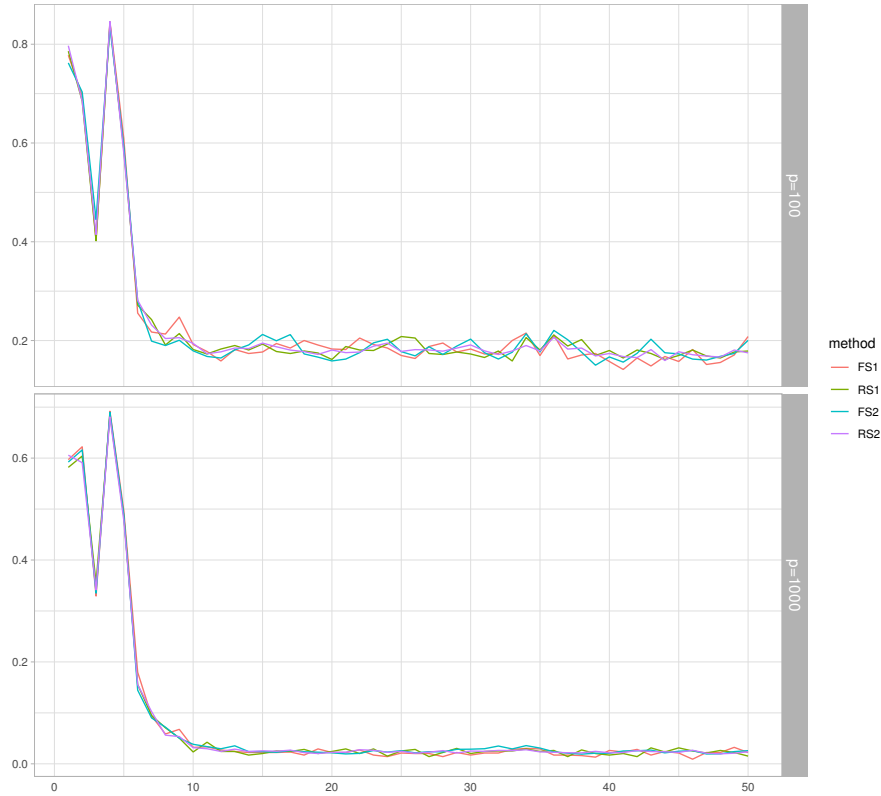


**Fig. S2.** ABC inclusion probabilities of the first 30 variables over different  $s$ . Each panel corresponds to a different combination of  $p \in \{100, 1000\}$  and  $M \in \{1000, 10000\}$ . Each row indicates a different model averaging strategy based on a different ratio of  $s$  over  $n$ . Each column corresponds to a different sample size. The legend represents various combinations of  $T \times B$ . For example,  $20 \times 200$  means each forest consists of  $T = 20$  trees and  $B = 200$  MCMC iterations as burnin. Note that we use  $\epsilon = \{\text{top } 10\%\}$  here.

mal training sample is very challenging and one needs to exercise statistical judgment to select from among various strategies. While Berger and Pericchi (2004) argue that: “Judgments involved in choosing good training samples will typically be much less than the judgments needed to implement an actual subjective Bayesian analysis”, we argue that entertaining some reasonable form of the data splitting (even if not optimal) will provide better results than naive ABC strategy in our context. The following simulated example shows that the variable selection performance with internal splitting is at least as good as with pre-determined splitting. We still use the Friedman’s dataset with



$n = 500$ ,  $p = 100$  and  $p = 1000$ ,  $\sigma = 5$  and autocorrelation 0.9. The ABC settings are  $T = 20$ ,  $\theta = 0.5$ ,  $s = 0.5n$ ,  $\epsilon = \text{top } 10\%$ . The inclusion probabilities are averaged over 10 datasets and plotted in Figure S3. From Figure S3, we can see that the differences in



**Fig. S3.** Comparison of Inclusion Probabilities of Pre-determined Splitting (FS) and Internal Splitting (RS). The inclusion probabilities are averaged over 10 independent Friedman’s datasets ( $n = 500$ ,  $\sigma = 5$ , autocorrelation = 0.9). FS1/RS1 are built with  $M = 1000$ , and FS2/RS2 are built with  $M = 10000$

inclusion probabilities of pre-determined splitting and internal splitting are small. This could be explained by the fact the data have been generated with Gaussian noise without any outliers which could potentially affect quality of splits.

Combining our findings from all of the sensitivity analyses above, we recommend the following default settings for the parameters:  $s = 0.5n$ ,  $T = 20$ , burnin = 200,  $\epsilon = \text{top } 10\%$ ,  $M = 1,000$  and internal data splitting.

### S.5. Full HIV Data Analysis

In this section, we provide a summary of our results on the entire dataset from Barber and Candès (2015). The summary statistics of the data are reported in Table S2. Comparisons are made between ABC Bayesian Forests, BART, DART and Random Forests. BART and DART are run with 50 trees for 20000 MCMC iterations (taking the first

**Table S2.** Basic summary statistics of the HIV dataset. DS refers to the decrease in susceptibility of the drug once the mutations has occurred.

HIV Virus Life Cycle	Drug Class	Mean Log DS	Number of Features	Number of Samples
PI	APV	0.75	201	767
	ATV	1.59	147	328
	IDV	1.33	206	825
	LPV	1.74	184	515
	NFV	2.00	207	842
	RTV	1.72	205	793
	SQV	1.22	206	824
NRTI	X3TC	3.10	283	629
	ABC	1.14	283	623
	AZT	1.55	283	626
	D4T	0.43	281	625
	DDI	0.43	283	628
	TDF	0.22	215	351
NNRTI	DLV	0.98	305	730
	EFV	1.08	312	732
	NVP	1.80	313	744

10 000 as a burn-in). Random Forests are implemented with the default number of 500 trees.

To summarize the results, we adopted 2 cutoff selection criteria. The first selection threshold is adaptive and is chosen as the maximum importance measure of a non-experimentally validated mutation. This cutoff point corresponds to zero false discoveries. Next, we use an automatic criterion for each method. For ABC Bayesian Forests (run with  $T = 20$  trees and  $M = 200$  burnin iterations, 10 000 ABC samples and top 100, 500 and 1 000 samples with the smallest discrepancy), we adopted the median probability model with the 0.5 cutoff. For DART and BART, we choose variables which have been split on at least once on average. For Random forest, the RFE approach (as described in [Linerio \(2018\)](#)) is used to find the variables. Similarly as in [Barber and Candès \(2015\)](#), we report the number true positions discovered and the number of false positions. To further study the separation power, we also report AUC of each method. The results are shown in [Table S3](#), [S4](#) and [S5](#).

Across all the drugs, we notice that ABC Bayesian Forest has a strong separation power, as is indicated by the performance of AUC scores. Random Forests with RFE tends to overfit by selecting too many mutations. BART and DART are performing well in this case but ABC is seen to have better AUC while being overall more conservative.

**Table S3.** The table summarizes results for a drug class PI. There are three performance criteria. For the adaptive cutoff, we report the number of true discoveries since the number of false discoveries is 0. For the automatic cutoff, we report both the number of false and true discoveries. Finally, we report a cutoff-free metric AUC. The best performance in each row is in bold font.

		APV					
Methods		ABC			BART	DART	Random Forest
		100	500	1000			
Adaptive cut-off	True Discoveries	17	<b>19</b>	<b>19</b>	14	15	15
Automatic cut-off	False Discoveries	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	7	31
	True Discoveries	13	11	11	14	20	<b>34</b>
AUC		0.69	0.75	<b>0.77</b>	0.65	0.65	0.61
		ATV					
Adaptive cut-off	True Discoveries	<b>23</b>	<b>23</b>	<b>23</b>	19	19	13
Automatic cut-off	False Discoveries	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	3	<b>0</b>
	True Discoveries	16	15	15	18	<b>21</b>	19
AUC		0.77	0.78	<b>0.79</b>	0.62	0.65	0.71
		IDV					
Adaptive cut-off	True Discoveries	8	9	9	6	11	<b>13</b>
Automatic cut-off	False Discoveries	<b>1</b>	<b>1</b>	<b>1</b>	2	5	32
	True Discoveries	14	14	14	18	18	<b>34</b>
AUC		0.73	<b>0.75</b>	<b>0.75</b>	0.65	0.63	0.62
		LPV					
Adaptive cut-off	True Discoveries	14	14	14	<b>15</b>	13	9
Automatic cut-off	False Discoveries	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	7	31
	True Discoveries	13	13	13	14	17	<b>34</b>
AUC		0.72	0.74	<b>0.75</b>	0.56	0.57	0.62
		NFV					
Adaptive cut-off	True Discoveries	8	10	10	11	<b>16</b>	15
Automatic cut-off	False Discoveries	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	5	32
	True Discoveries	15	15	14	17	20	<b>34</b>
AUC		0.73	<b>0.74</b>	<b>0.74</b>	0.65	0.64	0.65
		RTV					
Adaptive cut-off	True Discoveries	10	10	9	<b>13</b>	11	11
Automatic cut-off	False Discoveries	2	<b>1</b>	<b>1</b>	3	4	31
	True Discoveries	13	11	11	14	20	<b>34</b>
AUC		0.72	0.74	<b>0.75</b>	0.62	0.60	0.67
		SQV					
Adaptive cut-off	True Discoveries	15	15	15	3	<b>17</b>	10
Automatic cut-off	False Discoveries	<b>0</b>	<b>0</b>	<b>0</b>	3	6	31
	True Discoveries	15	15	14	16	17	<b>34</b>
AUC		0.74	0.77	<b>0.78</b>	0.64	0.62	0.57

**Table S4.** The table summarizes results for a drug class NRTI. There are three performance criteria. For the adaptive cutoff, we report the number of true discoveries since the number of false discoveries is 0. For the automatic cutoff, we report both the number of false and true discoveries. Finally, we report a cutoff-free metric AUC. The best performance in each row is in bold font.

		X3TC					
Methods		100	ABC		BART	DART	Random Forest
			500	1000			
Adaptive cut-off	True Discoveries	6	<b>9</b>	<b>9</b>	4	5	6
Automatic cut-off	False Discoveries	<b>0</b>	<b>0</b>	<b>0</b>	4	3	6
	True Discoveries	6	5	5	7	12	<b>15</b>
AUC		<b>0.70</b>	<b>0.70</b>	<b>0.70</b>	0.62	0.64	0.66
		ABC					
Adaptive cut-off	True Discoveries	8	8	7	7	10	<b>12</b>
Automatic cut-off	False Discoveries	2	<b>1</b>	1	<b>1</b>	7	2
	True Discoveries	10	10	10	11	14	<b>16</b>
AUC		0.74	0.73	<b>0.76</b>	0.66	0.71	0.74
		AZT					
Adaptive cut-off	True Discoveries	7	7	7	3	10	<b>13</b>
Automatic cut-off	False Discoveries	2	<b>1</b>	<b>1</b>	6	8	2
	True Discoveries	12	11	11	14	<b>16</b>	15
AUC		0.71	0.72	<b>0.73</b>	0.70	0.69	0.75
		D4T					
Adaptive cut-off	True Discoveries	<b>9</b>	8	<b>9</b>	5	0	8
Automatic cut-off	False Discoveries	2	<b>1</b>	<b>1</b>	3	12	80
	True Discoveries	12	12	11	12	14	<b>24</b>
AUC		<b>0.75</b>	<b>0.75</b>	<b>0.75</b>	0.70	0.70	0.73
		DDI					
Adaptive cut-off	True Discoveries	5	5	6	7	3	<b>10</b>
Automatic cut-off	False Discoveries	<b>1</b>	<b>1</b>	<b>1</b>	2	11	81
	True Discoveries	8	7	7	8	13	<b>24</b>
AUC		0.71	0.73	<b>0.74</b>	0.68	0.66	0.72
		TDF					
Adaptive cut-off	True Discoveries	4	<b>9</b>	<b>9</b>	3	7	2
Automatic cut-off	False Discoveries	2	<b>1</b>	<b>1</b>	4	11	8
	True Discoveries	9	9	9	10	18	<b>15</b>
AUC		0.69	0.72	0.72	0.72	<b>0.75</b>	0.73

**Table S5.** The table summarizes results for a drug class NNRTI. There are three performance criteria. For the adaptive cutoff, we report the number of true discoveries since the number of false discoveries is 0. For the automatic cutoff, we report both the number of false and true discoveries. Finally, we report a cutoff-free metric AUC. The best performance in each row is in bold font.

		DLV					
Methods		ABC			BART	DART	Random Forest
		100	500	1000			
Adaptive cut-off	True Discoveries	<b>4</b>	<b>4</b>	<b>4</b>	3	3	3
Automatic cut-off	False Discoveries	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	8	96
	True Discoveries	7	<b>7</b>	7	9	10	<b>14</b>
AUC		0.84	<b>0.87</b>	<b>0.87</b>	0.73	0.70	0.81
		EFV					
Adaptive cut-off	True Discoveries	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>	4	4
Automatic cut-off	False Discoveries	5	<b>4</b>	<b>4</b>	5	6	9
	True Discoveries	8	7	6	9	9	<b>10</b>
AUC		0.80	0.83	<b>0.84</b>	0.74	0.73	0.78
		NVP					
Adaptive cut-off	True Discoveries	6	6	6	8	6	<b>14</b>
Automatic cut-off	False Discoveries	3	3	<b>2</b>	<b>2</b>	9	97
	True Discoveries	6	6	5	<b>7</b>	6	5
AUC		0.79	0.79	0.79	0.71	0.66	<b>0.82</b>