
Generative Forests

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 We focus on generative AI for a type of data that still represent one of the most
2 prevalent form of data: tabular data. We introduce a new powerful class of forest-
3 based models fit for such tasks and a simple training algorithm with strong conver-
4 gence guarantees in a boosting model that parallels that of the original weak / strong
5 supervised learning setting. This algorithm can be implemented by a few tweaks
6 to the most popular induction scheme for decision tree induction (*i.e.* *supervised*
7 *learning*) with two classes. Experiments on the quality of generated data display
8 substantial improvements compared to the state of the art. The losses our algorithm
9 minimize and the structure of our models make them practical for related tasks that
10 require fast estimation of a density given a generative model and an observation
11 (even partially specified): such tasks include missing data imputation and density
12 estimation. Additional experiments on these tasks reveal that our models can be
13 notably good contenders to diverse state of the art methods, relying on models as
14 diverse as (or mixing elements of) trees, neural nets, kernels or graphical models.

15

1 Introduction

16 There is a substantial resurgence of interest in the ML community around tabular data, not just
17 because it is still one of the most prominent kind of data available [6]: it is recurrently a place of
18 sometimes heated debates on what are the best model architectures to solve related problems. For
19 example, even for well-posed problems with decades of theoretical formalization like supervised
20 learning [37], after more than a decade of deep learning disruption [21], there is still much ink spilled
21 in the debate decision forests vs neural nets [27]. Generative AI* makes no exception: where the
22 consensus has long been established on the best categories of architectures for data like image and
23 text (neural nets), tabular data still flourishes with a variety of model architectures building up – or
24 mixing – elements from knowledge representation, logics, kernel methods, graph theory, and of
25 course neural nets [4, 9, 12, 29, 32, 39, 41, 42] (see Section 2). Because of the remarkable nature of
26 tabular data where a single variable can bring considerable information about a target to model, each
27 of these classes can be a relevant choice at least in *some* cases (such is is the conclusion of [27] in the
28 context of supervised learning). A key differentiator between model classes is then training and the
29 formal guarantees it can provide [29, 42].

30 **In this paper**, we introduce new generative models based on sets of trees that we denote as *generative*
31 *forests* (GF), along with a training algorithm which has two remarkable features: it is extremely simple
32 and brings strong convergence guarantees in a weak / strong learning model that parallels that of the
33 original boosting model of Valiant’s PAC learning [17]. These guarantees improve upon the best state
34 of the art guarantees from several standpoint: [42] rely on convergence in the limit with restrictive
35 assumptions on the target density, [29] rely on a symmetry assumption on the loss optimized which
36 is restrictive for data generation. Our training algorithm, GF.BOOST, supports training from data

*We use this now common parlance expression on purpose, to avoid confusion with the other "generative" problem that consists in modelling densities [7, 35].

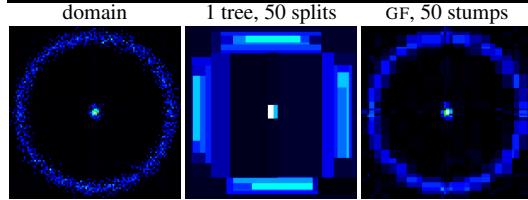


Table 1: *Left*: domain `circgauss`. *Center*: density learned by a generative forest (GF) consisting of a single tree, boosted for a small number (50) of iterations. *Right*: density learned by a GF consisting of 50 boosted tree stumps (*Center* and *Right* learned using GF.BOOTSTRAP). In a domain \mathcal{X} with dimension d , a single tree with n splits can only partition the domain in $n + 1$ parts. On the other hand, a set of n stumps in a GF can boost this number to $n^{\tilde{\Omega}(d)}$ (the tilda omits log dependencies), which explains why the right density appears so much better than the central one, even when each tree is just a stump.

37 with missing values and is simple enough to be implementable by a few tweaks on the popular
 38 induction scheme for *decision tree induction* with two classes, which is supported by a huge number
 39 of repositories / ML software implementing algorithms like CART or C4.5 [1, 33, 3, 34, 43]. From
 40 the model standpoint, generative forests bring a sizeable combinatorial advantage over its closest
 41 competitors, generative trees [29] (see Table 1 for an example) and adversarial random forests [42].
 42 Experiments on a variety of simulated or readily available domains display that our models can
 43 substantially improve upon state of the art. Our experiments also demonstrate that the observation of
 44 Table 1 can hold even for real-world domains as a small number of small trees (even stumps) in a
 45 generator can be enough to get state of the art results.
 46 The losses GF.BOOTSTRAP directly minimize are grounded in supervised learning and density estimation.
 47 The models we build have an additional benefit: we can efficiently compute the full density given
 48 an observation (*even partially specified*) and a generative model of ours. A GF can thus be used to
 49 generate data, but can also be used for side tasks like missing data imputation or density estimation.
 50 We provide additional experiments that clearly display that the same models of ours, trained by
 51 GF.BOOTSTRAP, can be notably good contenders to the state of the art.
 52 All proofs and additional experiments and results are given in an Appendix.

53 2 Related work

54 It would not do justice to the large amount of work in the field of "Generative AI" for tabular data
 55 to just sample a few of them, so we devote a part of the Appendix to an extensive review of the
 56 state of the art. Let us just mention that, unlike for unstructured data like images, there is a huge
 57 variety of model types, based on trees [7, 29, 42], neural networks [12, 14, 19, 45], probabilistic
 58 circuits [5, 41, 36], kernel methods [4, 32], graphical models [39] (among others: note that some
 59 are in fact hybrid models). The closest approaches to ours are [42] and [29], because the models
 60 include trees with a stochastic activation of edges to pick leaves, and a leaf-dependent data generation
 61 process. While [29] learn a single tree, [42] use a way to generate data from a set of trees – called an
 62 adversarial random forest – which is simple: sample a tree, and then sample an observation from the
 63 tree. The model is thus simple but not economical: each observation is generated by a single tree
 64 only, each of them thus having to represent a good sampler in its own. In our case, generating one
 65 observation makes use of *all* trees (Figure 1). We also note that the theory part of [42] is limited
 66 because it resorts to statistical consistency (infinite sample) and Lipschitz continuity of the target
 67 density, with second derivative continuous, square integrable and monotonic. Similarly, [29] resort to
 68 a wide set of proper losses, but with a symmetry constraint impeding in a generative context.

69 3 Definitions

70 Perhaps surprisingly at first glance, this Section introduces generative and *supervised* loss functions.
 71 Indeed, our algorithm, which trains a data generator and whose overall convergence shall be shown
 72 on a generative loss, turns out to locally optimize a *supervised* loss. For the interested reader, the key
 73 link, which is of independent interest since it links in our context losses for supervised and generative
 74 learning, is established in Lemma A (Appendix).
 75 $\forall k \in \mathbb{N}_{>0}$, let $[k] \doteq \{1, 2, \dots, k\}$. Our notations follow [35]. Let $\mathfrak{B} \doteq (\pi, A, B)$ denote a *binary*
 76 *task*, where A, B (and any other measure defined hereafter) are probability measures with the same
 77 support, also called *domain*, \mathcal{X} , and $\pi \in [0, 1]$ is a *prior*. $M \doteq \pi \cdot A + (1 - \pi) \cdot B$ is the corresponding

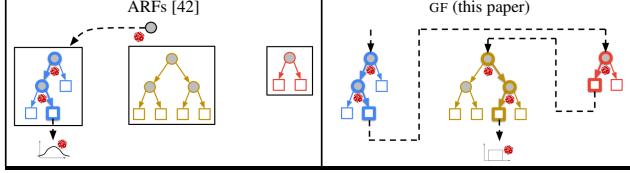


Figure 1: Sketch of comparison of two approaches to generate one observation, using Adversarial Random Forests [42] (left) and using generative forests, GF (right, this paper). In the case of Adversarial Random Forests, a tree is sampled uniformly at random, then a leaf is sampled in the tree and finally an observation is sampled according to the distribution "attached" to the leaf. Hence, only one tree is used to generate an observation. In our case, we leverage the combinatorial power of the trees in the forest: *all trees* are used to generate one observation, as each is contributing to one leaf. Figure 3 provides more details on generation using GF.

78 mixture measure. For the sake of simplicity, we assume \mathcal{X} bounded hereafter, and note that tricks can
 79 be used to remove this assumption [29, Remark 3.3]. In tabular data, each of the $\dim(\mathcal{X})$ features
 80 can be of various *types*, including categorical, numerical, etc., and associated to a natural measure
 81 (counting, Lebesgue, etc.) so we naturally associate \mathcal{X} to the product measure, which can thus be
 82 of mixed type. We also write $\mathcal{X} \doteq \times_{i=1}^d \mathcal{X}_i$, where \mathcal{X}_i is the set of values that can take on variable
 83 i . Several essential measures will be used in this paper, including \mathbf{U} , the uniform measure, \mathbf{G} , the
 84 measure associated to a generator that we learn, \mathbf{R} , the *empirical* measure corresponding to a training
 85 sample of observations. Like in [29], we do not investigate generalisation properties.

86 **Loss functions** There is a natural problem associated to binary task \mathfrak{B} , that of estimating the
 87 probability that an arbitrary observation $x \in \mathcal{X}$ was sampled from A – call such *positive* – or B – call
 88 these *negative* –. To learn a *supervised* model $\mathcal{X} \rightarrow [0, 1]$ for such a *class probability estimation*
 89 (CPE) problem, one usually has access to a set of examples where each is a couple (observation, class),
 90 the class being in set $\mathcal{Y} \doteq \{-1, 1\}$ ($=\{\text{negative, positive}\}$). Examples are drawn i.i.d. according
 91 to \mathfrak{B} . Learning a model is done by minimizing a loss function: when it comes to CPE, any such
 92 CPE loss [2] is some $\ell : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ whose expression can be split according to *partial* losses
 93 $\ell_1, \ell_{-1}, \ell(y, u) \doteq [\mathbb{y} = 1] \cdot \ell_1(u) + [\mathbb{y} = -1] \cdot \ell_{-1}(u)$. Its (pointwise) *Bayes risk* function is the
 94 best achievable loss when labels are drawn with a particular positive base-rate,

$$\underline{L}(p) \doteq \inf_u \mathbb{E}_{Y \sim B(p)} \ell(Y, u), \quad (1)$$

95 where $B \doteq$ Bernoulli random variable for class 1. A fundamental property for a CPE loss is *properness*,
 96 encouraging to guess ground truth: ℓ is proper iff $\underline{L}(p) = \mathbb{E}_{Y \sim B(p)} \ell(Y, p)$, $\forall p \in [0, 1]$, and *strictly*
 97 proper if $\underline{L}(p) < \mathbb{E}_{Y \sim B(p)} \ell(Y, u)$, $\forall u \neq p$. Strict properness implies strict concavity of Bayes risk.
 98 For example, the square loss has $\ell_1^{\text{sq}}(u) \doteq (1 - u)^2$, $\ell_{-1}^{\text{sq}}(u) \doteq u^2$, and, being strictly proper, Bayes
 99 risk $\underline{L}^{\text{sq}}(u) \doteq u(1 - u)$. Popular strictly proper ML include the log and Matusita's losses. All these
 100 losses are *symmetric* since $\ell_1^{\text{sq}}(u) = \ell_{-1}^{\text{sq}}(1 - u)$, $\forall u \in (0, 1)$ and *differentiable* because both partial
 101 losses are differentiable.

102 In addition to CPE losses, we introduce a set of losses relevant to generative approaches, that are
 103 popular in density ratio estimation [28, 38]. For any differentiable and convex $F : \mathbb{R} \rightarrow \mathbb{R}$, the
 104 Bregman divergence with generator F is $D_F(z \| z') \doteq F(z) - F(z') - (z - z')F'(z')$. Given function
 105 $g : \mathbb{R} \rightarrow \mathbb{R}$, the generalized perspective transform of F given g is $\check{F}(z) \doteq g(z) \cdot F(z/g(z))$, g being
 106 implicit in notation \check{F} [25, 26, 30]. The *Likelihood ratio risk* of B with respect to A for loss ℓ is

$$\mathbb{D}_\ell(A, B) \doteq \pi \cdot \mathbb{E}_{\mathbf{U}} \left[D_{(-\underline{L})} \left(\frac{dA}{d\mathbf{U}} \middle\| \frac{dB}{d\mathbf{U}} \right) \right], \quad (2)$$

107 with $g(z) \doteq z + (1 - \pi)/\pi$ in the generalized perspective transform. The prior multiplication is for
 108 technical convenience. \mathbb{D}_ℓ is non-negative; strict properness is necessary for a key property of \mathbb{D}_ℓ :
 109 $\mathbb{D}_{\underline{L}} = 0$ iff $A = B$ almost everywhere [29].

110 4 Generative forests: models and data generation

111 **Architecture** We first introduce the basic building block of our models, *trees*.

112 **Definition 4.1.** A *tree* Υ is a binary directed tree whose internal nodes are labeled with an observation
 113 variable and arcs are consistently labeled with subsets of their tail node's variable domain.

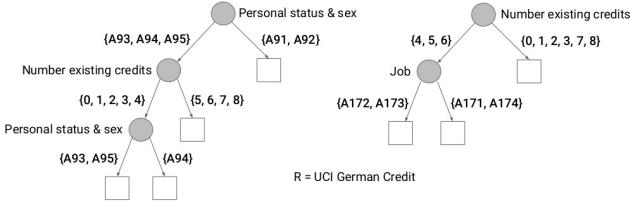


Figure 2: A GF ($T = 2$) associated to UCI German Credit.

Algorithm 1 INIT($\{\Upsilon_t\}_{t=1}^T$)

Input: Trees $\{\Upsilon_t\}_{t=1}^T$ of a GF;
Step 1 : **for** $t \in [T]$
 $\Upsilon_t.\nu^* \leftarrow \text{root}_t;$
 $\Upsilon_t.\text{done} \leftarrow [\Upsilon_t.\nu^* \in \Lambda(\Upsilon_t)];$

Algorithm 2 STARUPDATE($\Upsilon, \mathcal{C}, \mathbf{R}$)

Input: Tree Υ , subset $\mathcal{C} \subseteq \mathcal{X}$, measure \mathbf{R} ;
Step 1 : **if** $B(p_{\mathbf{R}} | \mathcal{X}_{\Upsilon.\nu^*} \cap \mathcal{C} | |\mathcal{C}|) = 1$ **then**
 $\mathcal{C} \leftarrow \mathcal{X}_{\Upsilon.\nu^*} \cap \mathcal{C}; \Upsilon.\nu^* \leftarrow \Upsilon.\nu_t^*;$
else
 $\mathcal{C} \leftarrow \mathcal{X}_{\Upsilon.\nu^*} \cap \mathcal{C}; \Upsilon.\nu^* \leftarrow \Upsilon.\nu_f^*;$
Step 2 : **if** $\Upsilon.\nu^* \in \Upsilon.\Lambda$ **then** $\Upsilon.\text{done} \leftarrow \text{true};$

114 *Consistency* is an important notion, linked to the generative model to which the tree belongs to.
115 We shall define it formally after having introduced our generative models. Informally however,
116 consistency postulates that the arcs' labels define a partition of the measure's support.

117 **Definition 4.2.** A *generative forest* (GF), \mathbf{G} , is a set of trees, $\{\Upsilon_t\}_{t=1}^T$, associated to measure \mathbf{R}
118 (implicit in notation).

119 Figure 2 shows an example of GF. We now proceed to defining the notion of consistency appearing in
120 Definition 4.1. For any node $\nu \in \mathcal{N}(\Upsilon)$ (the whole set of nodes of Υ , including leaves), we denote
121 $\mathcal{X}_{\nu} \subseteq \mathcal{X}$ the *support* of the node, which is \mathcal{X} for the root and then is computed as follows otherwise
122 for an internal node ν : starting with \mathcal{X} at the root, we progressively update the support as we descend
123 the tree, chopping off a feature's domain that is *not* labelling an arc followed, until we reach ν . Then,
124 a consistent labeling of arcs in a tree complies with one constraint:

125 **(C)** for each internal node ν and its left and right children ν_f, ν_t (respectively), $\mathcal{X}_{\nu} = \mathcal{X}_{\nu_f} \cup \mathcal{X}_{\nu_t}$
126 and the measure of $\mathcal{X}_{\nu_f} \cap \mathcal{X}_{\nu_t}$ with respect to \mathbf{G} is zero.

127 Hence, a single tree defines a recursive partition of \mathcal{X} according to the splits induced by the inner
128 nodes. Such is also the case for a set of trees, where *intersections* of the supports of tuples of leaves
129 (1 for each tree) define the subsets:

$$\mathcal{P}(\mathbf{G}) \doteq \{\cap_{i=1}^T \mathcal{X}_{\lambda_i} \text{ s.t. } \lambda_i \in \Lambda(\Upsilon_i), \forall i\} \quad (3)$$

130 $(\Lambda(\Upsilon) \subseteq \mathcal{N}(\Upsilon))$ is the set of leaves of Υ). Notably, we can construct the elements of $\mathcal{P}(\mathbf{G})$ using the
131 same algorithm that would compute it for 1 tree. First, considering a first tree Υ_1 , we compute the
132 support of a leaf, say \mathcal{X}_{λ_1} , using the algorithm described for the consistency property above. Then,
133 we start again with a second tree Υ_2 *but* replacing the initial \mathcal{X} by \mathcal{X}_{λ_1} , yielding $\mathcal{X}_{\lambda_1} \cap \mathcal{X}_{\lambda_2}$. Then
134 we repeat with a third tree Υ_3 replacing \mathcal{X} by $\mathcal{X}_{\lambda_1} \cap \mathcal{X}_{\lambda_2}$, and so on until the last tree is processed.
135 This yields one element of $\mathcal{P}(\mathbf{G})$.

136 **Generating one observation** Generating one observation relies on a *stochastic* version of the
137 procedure just described. It ends up in an element of $\mathcal{P}(\mathbf{G})$ of positive measure, from which we
138 sample uniformly one observation, and then repeat the process for another observation. To describe
139 the process at length, we make use of two key routines, INIT and STARUPDATE, see Algorithms
140 1 and 2. INIT initializes "special" nodes in each tree, that are called *star nodes*, to the root of
141 each tree (notation for a variable v relative to a tree Υ is $\Upsilon.v$). Stochastic activation, performed in
142 STARUPDATE, progressively makes star nodes descend in trees. When all star nodes have reached a
143 leaf in their respective tree, an observation is sampled from the intersection of the leaves' domains
144 (which is an element of $\mathcal{P}(\mathbf{G})$). A Boolean flag, done takes value *true* when the star node is in the
145 leaf set of the tree ($[\cdot]$ is Iverson's bracket for the truth value of a predicate, [20]).

146 The crux of generation is STARUPDATE. This procedure is called with a tree of the GF *for which*
147 *done is false*, a subset \mathcal{C} of the whole domain and measure \mathbf{R} . The first call of this procedure is done
148 with $\mathcal{C} = \mathcal{X}$. When all trees are marked done, \mathcal{C} has been "reduced" to some $\mathcal{C}_s \in \mathcal{P}(\mathbf{G})$, where the
149 index reminds that this is the last \mathcal{C} we obtain, from which we sample an observation, uniformly at
150 random in \mathcal{C}_s . Step 1 in STARUPDATE is fundamental: it relies on tossing an unfair coin (a Bernoulli

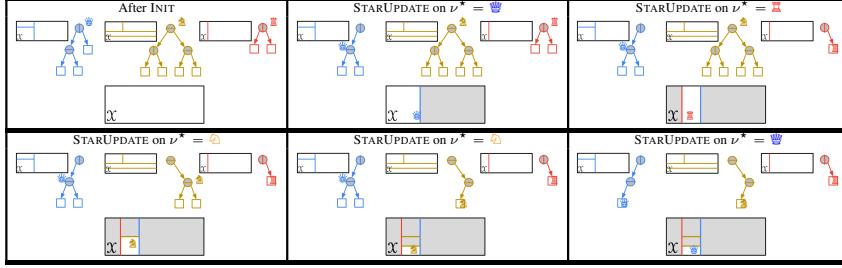


Figure 3: From left to right and top to bottom: updates of the argument \mathcal{C} of STARUPDATE through a sequence of run of STARUPDATE in a generative forest consisting of three trees (the partition of the domain induced by each tree is also depicted, alongside the nature of splits, vertical or horizontal, at each internal node) whose star nodes are indicated with chess pieces (\spadesuit , \clubsuit , \heartsuit). In each picture, \mathcal{C} is represented at the bottom of the picture (hence, $\mathcal{C} = \mathcal{X}$ after INIT). In the bottom-right picture, all star nodes are leaves and thus $\mathcal{C}_s = \mathcal{C}$ displays the portion of the domain in which an observation is sampled. Remark that the last star node update produced no change in \mathcal{C} . The *sequence* of choice of trees does not matter as the key invariant of Lemma 4.4 is always satisfied through STARUPDATE: the probability to get \mathcal{C} is always equal to its empirical mass in \mathbf{R} . The Appendix (Section II) details the sequential, randomized and concurrent choices. On this Figure, the choice is not sequential.

151 event noted $B(p)$), where the head probability p is just the mass of \mathbf{R} in $\mathcal{X}_{\mathbf{r}, \nu_t^*} \cap \mathcal{C}$ relative to \mathcal{C} .
 152 Hence, if $\mathcal{X}_{\mathbf{r}, \nu_t^*} \cap \mathcal{C} = \mathcal{C}$, $p = 1$. There is a simple but important invariant (proof omitted).

153 **Lemma 4.3.** *In STARUPDATE, it always holds that the input \mathcal{C} satisfies $\mathcal{C} \subseteq \mathcal{X}_{\mathbf{r}, \nu^*}$.*

154 Remark that we have made no comment about the *sequence* of tree choices over which STARUP-
 155 DATE is called. Let us call *admissible* such a sequence that ends up with *some* $\mathcal{C}_s \subseteq \mathcal{X}$. T being the
 156 number of trees (see INIT), for any sequence $v \in [T]^{D(\mathcal{C}_s)}$, where $D(\mathcal{C}_s)$ is the sum of the depths of
 157 all the star leaves whose support intersection is \mathcal{C}_s , we say that v is *admissible for* \mathcal{C}_s if there exists a
 158 sequence of branchings in Step 1 of STARUPDATE, whose corresponding sequence of trees follows
 159 the indexes in v , such that at the end of the sequence all trees are marked done and the last $\mathcal{C} = \mathcal{C}_s$.
 160 We note that if we permute two different indexes in v , the sequence is still admissible for the same \mathcal{C}_s
 161 but if we *flip* an index for another one, it is not admissible anymore for \mathcal{C}_s . Crucially, the probability
 162 to end up in \mathcal{C}_s using STARUPDATE, given *any* of its admissible sequences, is the *same* and equals its
 163 mass with respect to \mathbf{R} .

164 **Lemma 4.4.** *For any $\mathcal{C}_s \in \mathcal{P}(\mathbf{G})$ and admissible sequence $v \in [T]^{D(\mathcal{C}_s)}$ for \mathcal{C}_s , $p_{\mathbf{G}}[\mathcal{C}_s | v] = p_{\mathbf{R}}[\mathcal{C}_s]$.*
 165 The Lemma is simple to prove but fundamental in our context as the way one computes the sequence
 166 – and thus the way one picks the trees – does not bias generation: the sequence of tree choices could
 167 thus be iterative, randomized, concurrent (e.g. if trees were distributed), etc., this would not change
 168 generation’s properties from the standpoint of Lemma 4.4. We defer to Appendix (Section II) three
 169 examples of sequence choice. Figure 3 illustrates a sequence and the resulting \mathcal{C}_s .

170 **Missing data imputation and density estimation with GFS** A generative forest is not just a generator:
 171 it models a density and the exact value of this density at any observation is easily available.
 172 Figure 4 (top row) shows how to compute it. Note that if carried out in parallel, the complexity
 173 to get this density is cheap, of order $O(\max_t \text{depth}(\Upsilon_t))$. If one wants to prevent predicting zero
 174 density, one can stop the descent if the next step creates a support with zero empirical measure. A
 175 GF thus also models an easily tractable density, but this fact alone is not enough to make it a *good*
 176 density estimator. To get there, one has to factor the loss minimized during training. In our case,
 177 as we shall see in the following Section, we train a GF by minimizing an information-theoretic loss
 178 directly formulated on this density (2). So, using a GF also for density estimation can be a reasonable
 179 additional benefit of training GFs.

180 The process described above that finds leaves reached by an observation can be extended to missing
 181 values in the observation using standard procedures for classification using decision trees. We obtain
 182 a simple procedure to carry out *missing data imputation* instead of density estimation: once a tuple of
 183 leaves is reached, one uniformly samples the missing features in the corresponding support. This is
 184 fast, but at the expense of a bit more computations, we can have a much better procedure, as explained
 185 in Figure 4 (bottom row). In a first step, we compute the density of each support subset corresponding
 186 to *all* (not just 1 as for density estimation) tuples of leaves reached in each tree. This provides us
 187 with the *full density* over the missing values *given* (i) a partial assignment of the tabular domain’s
 188 variables and (ii) the GF. We then keep the elements corresponding to the maximal density value and

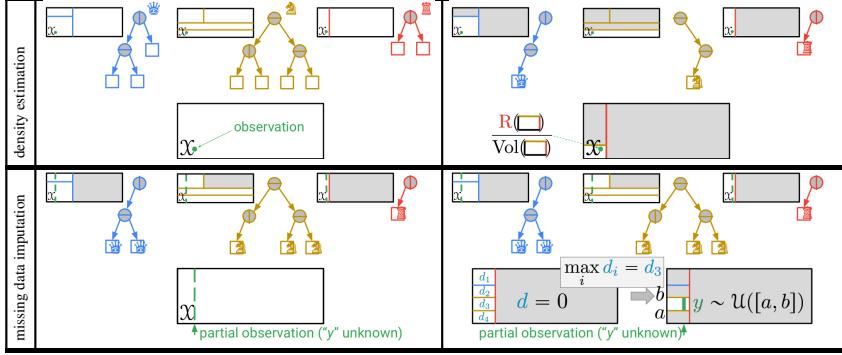


Figure 4: **(Top row)** Density estimation using a GF, on an observation indicated by \bullet (*Left*). In each tree, the leaf reached by the observation is found and the intersection of all leaves' supports is computed. The estimated density at \bullet is computed as the empirical measure in this intersection over its volume (*Right*). **(Bottom row)** Missing data imputation using the same GF, and an observation with one missing value (if $X \subset \mathbb{R}^2$, then y is missing). We first proceed like in density estimation, finding in each tree *all* leaves *potentially* reached by the observation if y were known (*Left*); then, we compute the density in *each* non-empty intersection of all leaves' supports; among the corresponding elements with maximal density, we get the missing value(s) by uniform sampling (*Right*).

Algorithm 3 GF.BOOST(\mathbf{R}, J, T)

```

Input: measure  $\mathbf{R}$ , #iters  $J$ , #trees  $T$ ;
Output: trees  $\{\Upsilon_t\}_{t=1}^T$  of GF  $\mathbf{G}$ ;
Step 1 :  $\mathcal{T} \leftarrow \{\Upsilon_t = (\text{root})\}_{t=1}^T$ ;
Step 2 : for  $j = 1$  to  $J$ 
    Step 2.1 :  $\Upsilon_* \leftarrow \text{tree } (\mathcal{T})$ ;
    Step 2.2 :  $\lambda_* \leftarrow \text{leaf } (\Upsilon_*)$ ;
    Step 2.3 :  $p \leftarrow \text{splitPred } (\lambda_*, \mathcal{T}, \mathbf{R})$ ;
    Step 2.4 :  $\text{split}(\Upsilon_*, \lambda_*, p)$ ;
return  $\mathcal{T}$ ;

```

189 finally simultaneously sample all missing features uniformly in the corresponding domain. Overall,
190 the whole procedure is $O(d \cdot (\sum_t \text{depth}(\Upsilon_t))^2)$.
191

192 5 Learning generative forests: boosting

193 In their paper, [29] show how to train a generative tree with two algorithms in the GAN framework,
194 requiring the help of a discriminator and repeated sampling using the generator. The way we train
195 our generative models provides two improvements to their setting: (i) we do not require sampling
196 the generator during training anymore, and more importantly (ii) we get rid of the discriminator,
197 thus training only the generator. Our training setting is not anymore generative as for GANs, but
198 *supervised* as our generative models are trained to minimize loss functions for a task called *class*
199 *probability estimation* [35], whose key property, *properness* [37], helped design all decades-old major
200 top-down decision tree induction algorithms [1, 34].

201 To learn a GF, we just have to learn its set of trees. Our training algorithm, GF.BOOST (Algorithm 3),
202 performs a greedy top-down induction. In Step 1, we initialize the set of T trees to T roots. Step
203 2.2 picks a candidate leaf to split in the chosen tree Υ_* and Step 2.4 splits Υ_* by replacing λ_* by a
204 stump whose corresponding splitting predicate, p , is returned in Step 2.3 using a weak splitter oracle
205 called *splitPred*. "weak" refers to boosting's weak/strong learning setting [17] and means that we
206 shall only require lightweight assumptions about this oracle; in decision tree induction, this oracle is
207 the key to boosting from such weak assumptions [18]. This will also be the case for our generative
208 models. We now investigate Step 2.3 and *splitPred*.

209 **The weak splitter oracle *splitPred*** In decision tree induction, a splitting predicate is chosen
210 to reduce an expected Bayes risk (1) (*e.g.* that of the log loss [34], square loss [1], etc.). In
211 our case, *splitPred* does about the same *with a catch in the binary task it addresses*, which is[†]
212 $\mathfrak{B}_{\text{GEN}} \doteq (\pi, \mathbf{R}, \mathbf{U})$ (our mixture measure is thus $\mathbf{M} \doteq \pi \cdot \mathbf{R} + (1 - \pi) \cdot \mathbf{U}$). The corresponding

[†]The prior is chosen by the user: without reasons to do otherwise, a balanced approach suggests $\pi = 0.5$.

213 expected Bayes risk that `splitPred` seeks to minimize is just:

$$\underline{\mathbb{L}}(\mathcal{T}) \doteq \sum_{\mathcal{C} \in \mathcal{P}(\mathcal{T})} p_{\text{M}}[\mathcal{C}] \cdot \underline{L}\left(\frac{\pi p_{\text{R}}[\mathcal{C}]}{p_{\text{M}}[\mathcal{C}]}\right). \quad (4)$$

214 The concavity of \underline{L} implies $\underline{\mathbb{L}}(\mathcal{T}) \leq \underline{\mathbb{L}}(\pi)$. Notation $\mathcal{P}(\cdot)$ overloads that in (3) by depending explicitly
 215 on the set of trees of $\textcolor{blue}{G}$ instead of G itself. In classical decision tree induction, (4) is optimized over a
 216 single tree: there is thus a single element in $\mathcal{P}(\mathcal{T})$ which is affected by the split, \mathcal{X}_{λ_*} , which makes the
 217 optimization of $\underline{\mathbb{L}}(\mathcal{T})$ very efficient from an implementation standpoint. In our case, multiple elements
 218 in $\mathcal{P}(\mathcal{T})$ can be affected by one split, so the optimisation is no more computationally demanding.
 219 From the standpoint of the potential decrease of $\underline{\mathbb{L}}(\mathcal{T})$ however, a single split in our case can be much
 220 more impactful. To see this, remark that for the candidate leaf λ_* ,

$$\sum_{\substack{\mathcal{C} \in \mathcal{P}(\mathcal{T}) \\ \mathcal{C} \subseteq \mathcal{X}_{\lambda_*}}} p_{\text{M}}[\mathcal{C}] \cdot \underline{L}\left(\frac{\pi p_{\text{R}}[\mathcal{C}]}{p_{\text{M}}[\mathcal{C}]}\right) = p_{\text{M}}[\mathcal{X}_{\lambda_*}] \sum_{\substack{\mathcal{C} \in \mathcal{P}(\mathcal{T}) \\ \mathcal{C} \subseteq \mathcal{X}_{\lambda_*}}} \frac{p_{\text{M}}[\mathcal{C}]}{p_{\text{M}}[\mathcal{X}_{\lambda_*}]} \cdot \underline{L}\left(\frac{\pi p_{\text{R}}[\mathcal{C}]}{p_{\text{M}}[\mathcal{C}]}\right) \leq p_{\text{M}}[\mathcal{X}_{\lambda_*}] \cdot \underline{L}\left(\frac{\pi p_{\text{R}}[\mathcal{X}_{\lambda_*}]}{p_{\text{M}}[\mathcal{X}_{\lambda_*}]}\right)$$

221 (because \underline{L} is concave). The leftmost term is the contribution of λ_* to $\underline{\mathbb{L}}(\mathcal{T})$, the rightmost its
 222 contribution to $\underline{\mathbb{L}}(\{\Upsilon_*\})$. In the latter case, the split gets two new terms instead of one, whose
 223 expectation is no larger than the contribution of λ_* to $\underline{\mathbb{L}}(\{\Upsilon_*\})$. In the former case, there can be up
 224 to $2 \cdot \text{Card}(\{\mathcal{C} \in \mathcal{P}(\mathcal{T}) : \mathcal{C} \subseteq \mathcal{X}_{\lambda_*}\})$ new contributions, so the inequality can have a substantial slack.
 225 **Boosting** Two questions remain: can we quantify the slack in decrease and of course what quality
 226 guarantee does it bring for the *generative* model $\textcolor{blue}{G}$ whose set of trees is learned by GF.BOOST ? We
 227 answer both questions in a single Theorem, which relies on an assumption that parallels the classical
 228 weak learning assumption of boosting:

229 **Definition 5.1. (WLA(γ, κ))** There exists $\gamma > 0, \kappa > 0$ such that at each iteration of GF.BOOST,
 230 the couple (λ_*, p) chosen in Steps 2.2, 2.3 of GF.BOOST satisfies the following properties: (a) λ_* is
 231 not skinny: $p_{\text{M}}[\mathcal{X}_{\lambda_*}] \geq 1/\text{Card}(\Lambda(\Upsilon_*))$, (b) truth values of p moderately correlates with $\mathfrak{B}_{\text{GEN}}$ at
 232 λ_* : $|p_{\text{R}}[p|\mathcal{X}_{\lambda_*}] - p_{\text{U}}[p|\mathcal{X}_{\lambda_*}]| \geq \gamma$, and finally (c) there is a minimal proportion of real data at λ_* :
 233 $\pi p_{\text{R}}[\mathcal{X}_{\lambda_*}]/p_{\text{M}}[\mathcal{X}_{\lambda_*}] \geq \kappa$.

234 The convergence proof of [18] reveals that both (a) and (b) are jointly made at any split, where our
 235 (b) is equivalent to their *weak hypothesis assumption* where their γ parameter is twice ours. (c)
 236 makes sense with our choice of $\mathfrak{B}_{\text{GEN}}$: the uniform measure of any node with non-empty support is
 237 always strictly positive, so we just guarantee that we do not split a node that would be useless for data
 238 generation. If ℓ is strictly proper and differentiable, then $\exists \kappa : \inf\{\ell'_{-1} - \ell'_1\} \geq \kappa > 0$ [24, Remark 1].
 239

240 **Theorem 5.2.** Suppose the loss ℓ is **strictly proper** and **differentiable**. Let $\textcolor{blue}{G}_0$ denote the initial
 241 GF with T roots in its trees (Step 1, GF.BOOST) and $\textcolor{blue}{G}_J$ the final GF, assuming wlog that the number
 242 of boosting iterations J is a multiple of T . Under WLA(γ, κ), we get the following on likelihood ratio
 243 risks: $\mathbb{D}_\ell(\mathbf{R}, \textcolor{blue}{G}_J) \leq \mathbb{D}_\ell(\mathbf{R}, \textcolor{blue}{G}_0) - \frac{\kappa \gamma^2 \kappa^2}{8} \cdot T \log(1 + \frac{J}{T})$.

244 Strict properness is essential for the loss to behave properly (Section 3) so only differentiability could
 245 be alleviated to generalize further our result, noting that popular non-differentiable losses are proper
 246 but not strictly proper (e.g. the 0/1 loss).

247 6 Experiments

248 Our code is provided and commented in Appendix, Section V.2. The main setting of our experiments
 249 is realistic data generation ('LIFELIKE'), but we have also tested our method for missing data
 250 imputation ('IMPUTE') and density estimation ('DENSITY'): for this reason, we have selected a
 251 broad panel of state of the art approaches to test against, relying on models as diverse as (or mixing
 252 elements of) trees, neural nets, kernels or graphical models, with MICE [40], adversarial random
 253 forests (ARFs [42]), CT-GANs [45], Vine copulas auto-encoders (VCAE, [39]) and Kernel density
 254 estimation (KDE, [4, 32]). We carried out experiments on a total of 21 datasets, from UCI [10],
 255 Kaggle, OpenML, the Stanford Open Policing Project, or simulated. All are presented in Appendix,
 256 Section V.1. We summarize results that are presented *in extenso* in Appendix. Before starting, we
 257 complete the 2D heatmap of Table 1 by another one showing our models can also learn deterministic
 258 dependences in real-world domains (Table 3). The Appendix also provides an example experiment
 259 on interpreting our models for a sensitive domain (in Section V.V.3.1).

domain tag	us (GF) $T=500, J=2\,000$	Adversarial Random Forests (ARFs)							
		$T = 10$	$p\text{-val}$	$T = 50$	$p\text{-val}$	$T = 100$	$p\text{-val}$	$T = 200$	$p\text{-val}$
iris	0.423 \pm 0.033	0.530 \pm 0.064	0.0039	0.517 \pm 0.081	0.0866	0.466 \pm 0.050	0.1452	0.502 \pm 0.043	0.0044
ring	0.285 \pm 0.008	0.289 \pm 0.006	0.1918	0.286 \pm 0.007	0.8542	0.288 \pm 0.010	0.7205	0.286 \pm 0.007	0.6213
circ	0.351 \pm 0.005	0.354 \pm 0.002	0.2942	0.356 \pm 0.008	0.2063	0.350 \pm 0.002	0.9050	0.355 \pm 0.005	0.3304
grid	0.390 \pm 0.002	0.394 \pm 0.003	0.1376	0.391 \pm 0.002	0.3210	0.392 \pm 0.001	0.2118	0.394 \pm 0.002	0.0252
for	1.108 \pm 0.105	1.272 \pm 0.281	0.2807	1.431 \pm 0.337	0.1215	1.311 \pm 0.255	0.0972	1.268 \pm 0.209	0.0616
rand	0.286 \pm 0.003	0.286 \pm 0.003	0.9468	0.290 \pm 0.005	0.2359	0.289 \pm 0.005	0.3990	0.288 \pm 0.002	0.3685
tic	0.575 \pm 0.002	0.577 \pm 0.001	0.2133	0.578 \pm 0.001	0.1104	0.577 \pm 0.001	0.2739	0.577 \pm 0.001	0.2773
iono	1.167 \pm 0.074	1.431 \pm 0.043	0.0005	1.469 \pm 0.079	0.0001	1.431 \pm 0.058	0.0015	1.420 \pm 0.056	0.0035
stm	0.975 \pm 0.026	1.010 \pm 0.015	0.0304	1.015 \pm 0.018	0.0231	1.001 \pm 0.010	0.0956	1.014 \pm 0.021	0.0261
wred	0.980 \pm 0.032	1.086 \pm 0.036	0.0003	1.072 \pm 0.055	0.0133	1.086 \pm 0.030	0.0003	1.099 \pm 0.031	0.0001
stp	0.976 \pm 0.018	0.999 \pm 0.011	0.0382	1.012 \pm 0.010	0.0250	1.006 \pm 0.011	0.0280	1.003 \pm 0.003	0.0277
ana	0.368 \pm 0.014	0.370 \pm 0.012	0.7951	0.382 \pm 0.015	0.1822	0.368 \pm 0.006	0.9161	0.369 \pm 0.007	0.8316
aba	0.490 \pm 0.028	0.505 \pm 0.046	0.5251	0.484 \pm 0.024	0.7011	0.503 \pm 0.035	0.6757	0.481 \pm 0.031	0.4632
wwhi	1.064 \pm 0.003	1.159 \pm 0.011	ε	1.159 \pm 0.006	ε	1.170 \pm 0.014	ε	1.150 \pm 0.021	0.0005
comp	0.532 \pm 0.008	0.531 \pm 0.012	0.8278	0.534 \pm 0.009	0.8438	0.551 \pm 0.023	0.1371	0.535 \pm 0.033	0.8642
arti	0.821 \pm 0.004	0.849 \pm 0.004	0.0007	0.847 \pm 0.010	0.0005	0.843 \pm 0.009	0.0006	0.849 \pm 0.010	0.0005
jung	0.929 \pm 0.002	0.929 \pm 0.002	0.6044	0.930 \pm 0.002	0.6536	0.929 \pm 0.002	0.5937	0.930 \pm 0.001	0.2447
elec	1.483 \pm 0.033	1.577 \pm 0.088	0.1052	1.543 \pm 0.017	0.0081	1.552 \pm 0.039	0.0522	1.581 \pm 0.025	0.0095
wins / lose for us →		15 / 1		17 / 1		16 / 1		17 / 1	

Table 2: LIFELIKE: comparison of ARFs [42] with different number T of trees, and Generative Forests (us, GF) where the number of trees and number of iterations are **fixed** ($T = 500$ trees, total of $J = 2000$ splits in trees). Values shown = average over the 5-folds \pm std dev. . The p -values for the comparison of our results ("us") and ARFs are shown. Values appearing in green in p -vals mean (i) the average regularised OT cost for us (GF) is smaller than that of ARFs and (ii) $p < 0.1$, i.e. our method outperforms ARFs (there is no instance of us being statistically significantly beaten by ARFs). ε means $p\text{-val} < 10^{-4}$. Domain details in Appendix, Table A1.

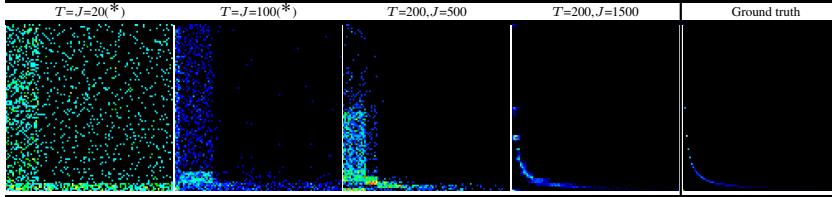


Table 3: 2D density plots generated on OpenML kc1 ($x = d$ and $y = l$) with GF, for varying number of trees T and number of splits J (columns). "*" = all trees are stumps. The rightmost column recalls the domain ground truth for comparison. Each generated dataset contains $m = 2000$ observations. The two variables (among $d = 22$) were chosen because of their deterministic dependence.

260 **Generation capabilities of our models: LIFELIKE** The objective of the experiment is to evaluate
 261 whether a generative model is able to create "realistic" data. The evaluation pipeline is simple: we
 262 create for each domain a 5-fold stratified experiment. Evaluating generation implies comparing
 263 distributions so we use as metric an optimal transport (OT) distance between the generated sample
 264 and the fold's test sample [8], constraining the generated data to be of the same size as the test fold.
 265 We compare our method to three possible contenders: Adversarial Random Forests (ARF) [42],
 266 CT-GANs [45] and Vine copula autoencoders (VCAE) [39]. All these approaches rely on models that
 267 are very different from each other. ARFs are trained with a number of trees $T \in \{10, 50, 100, 200\}$.
 268 ARFs include an algorithm to select the tree size so we do not have to select it. CT-GANs are trained
 269 with a number of epochs $E \in \{10, 100, 300, 1000\}$. We tested VCAE with all types of vine copulas
 270 (center, direct and regular). We compare those models with generative forests with $T = 500$ trees
 271 and trained for a total of $J = 2000$ iterations. We provide here the Table of our results vs the best
 272 contender (ARFs) and summarize the Tables vs the two others (given *in extenso* in Appendix).
 273 **Results vs Adversarial Random Forests** See Table 2. Globally, we manage to beat ARFs on almost
 274 all runs, and very significantly on many of them (we are never statistically beaten by ARFs). The
 275 Appendix demonstrates that smaller GFs can still be competitive though there is a domain-dependent
 276 minimal size for competitiveness. "Guessing" the right size is a prospect for future work.
 277 **Results summary vs CT-GANs and VCAE** Table 4 (left and center) aggregates all results vs CT-GANs
 278 and VCAE (see Appendix; parameters for GF are the same as in Table 2). We see that GF clearly beat
 279 those two approaches. Just like for ARFs, we are never statistically beaten by any of these methods.
 280 **Missing data imputation with our models: IMPUTE** We compared our method against a powerful
 281 contender, MICE [40], which relies on round-robin prediction of missing values using supervised
 282 models. We optimized MICE by choosing as supervised models trees (CART) and random forests (RFs,

GF >> CT-GAN	neither	CT-GAN >> GF	GF >> VCAE	neither	VCAE >> GF	GF >> KDE	neither	KDE >> GF
54	2	0	34	1	0	9	5	3

Table 4: *Left and Center table*: LIFELIKE, summary of the results vs CT-GANs and VCAE (Full Tables in Appendix: A6 and A7). We count the number of times we are statistically (significance explained in Table 2) better (left, these would be green cells in Table 2) or worse (right) than a contender. The central column counts the remaining times for which no statistical significance holds. *Right table*: DENSITY: comparison between us and KDE (summary), counting the number of domains for which we are statistically better (left), or worse (right). The central column counts the remaining domains, for which no statistical significance ever holds.

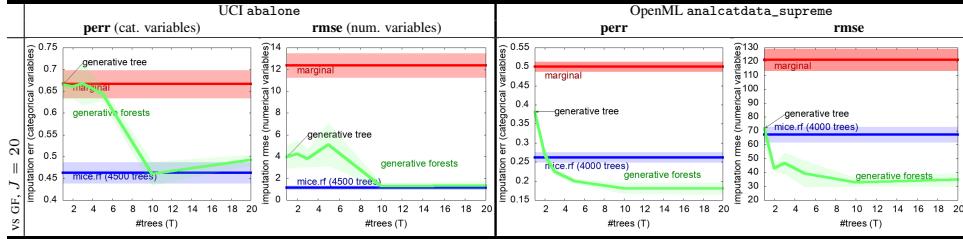


Table 5: IMPUTE: results on two domains (left and right pane). Since both domain contain categorical and numerical variables, we compute for each **perr** (categorical variables) and the **rmse** (numerical variables). In each plot, we display both the corresponding results of GF.BOOST, the results of MICE.RF (indicating the total number of trees used to impute the dataset) and the result of a fast imputation baseline using the **marginals** computed from the training sample. The *x* axis displays the number of trees T in GF.BOOST and each curve is displayed with its average \pm standard deviation in shaded color. The result for $T = 1$ equivalently indicates the performance of a single generative tree (GT) with J splits [29], shown with an arrow (see text).

we increased the number of trees to 100 for better results). We create a 5-fold experiment; in each fold, we remove a fixed % of observed values (data missing completely at random, MCAR). We then use the data *with missing values* as input to MICE or us to learn models which are then used to predict the missing values. We compute a per-feature discrepancy, the average error probability (**perr**, for categorical variables), and the root mean square error (**rmse**, numerical variables). We also compute one of the simplest baselines, which consists in imputing each variable from its empirical **marginal**. *Results summary* Table 5 puts the spotlight on two domains. The Appendix provides many more plots. From all results, we conclude that generative forests can compete or beat MICE.RF while using hundred times less trees (eventually using just stumps, when $J = T = 20$). From all our results, we also conclude that there is a risk of overfitting if J and / or T are too large. This could motivate further work on pruning generative models. For GF, an unexpected pattern is that the pattern "small number of small trees works well" can work on real-world domains as well (see also Appendix), thereby demonstrating that the nice result displayed in Table 1 generalises to bigger domains.

Density estimation: DENSITY We compared GF vs kernel density estimation (KDE) [4, 23, 32]. The experimental setting is the same as for LIFELIKE: in each of the 5-fold stratified cross validation fold, we use the training sample to learn a model (with KDE or GF.BOOST) which is then used to predict the observation's density in the the test sample. The higher the density, the better the model. The GF models we consider are the same as in the LIFELIKE ($T = 500$, $J = 2000$).

Results summary Table (4, right) summarizes the results (plots: Appendix, Section V.V.3.7). The leftmost column ("GF >> KDE") counts domains for which there exists an iteration j for GF after which we tend to be statistically better (and never statistically worse) than KDE. The rightmost column ("KDE >> GF") counts domains for which KDE is always statistically better than GF. The last type of plots appears to be those for which there is no such statistical difference (central column). The results support the conclusion that GF can also be good models to carry out density estimation.

7 Conclusion

We propose new models and a formal boosting algorithm for generative models based on ensemble of trees, with experiments on data generation, missing data imputation and density estimation. A responsible use of such models necessarily includes restriction to tabular data (Section 1) and one important related question, left for future work, is how to "guess" the right model size other than with cross-validation. A pruning algorithm with good generalization seems a very promising direction.

313 **References**

- 314 [1] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and regression trees*.
315 Wadsworth, 1984.
- 316 [2] A. Buja, W. Stuetzle, and Y. Shen. Loss functions for binary class probability estimation and
317 classification: structure and applications, 2005. Technical Report, University of Pennsylvania.
- 318 [3] T. Chen and C. Guestrin. XGBoost: A scalable tree boosting system. In *22nd KDD*, pages
319 785–794, 2016.
- 320 [4] Y.-C. Chen. A tutorial on kernel density estimation and recent advances. *Biostatistics &*
321 *Epidemiology*, 1:161–187, 2017.
- 322 [5] Y. Choi, A. Vergari, and G. Van den Broeck. Probabilistic circuits: a unifying framework for
323 tractable probabilistic models, 2020. <http://starai.cs.ucla.edu/papers/ProbCirc20.pdf>.
- 325 [6] M. Chui, J. Manyika, M. Miremadi, N. Henke, R. Chung P. Nel, and S. Malhotra. *Notes from*
326 *the AI frontier*. McKinsey Global Institute, 2018.
- 327 [7] A.-H.-C. Correia, R. Peharz, and C.-P. de Campos. Joints in random forests. In *NeurIPS’20*,
328 2020.
- 329 [8] M. Cuturi. Sinkhorn distances: lightspeed computation of optimal transport. In *NIPS’26*, pages
330 2292–2300, 2013.
- 331 [9] A. Darwiche. A logical approach to factoring belief networks. In Dieter Fensel, Fausto
332 Giunchiglia, Deborah L. McGuinness, and Mary-Anne Williams, editors, *KR’02*, pages 409–
333 420. Morgan Kaufmann, 2002.
- 334 [10] D. Dua and C. Graff. UCI machine learning repository, 2021.
- 335 [11] V. Dumoulin, I. Belghazi, B. Poole, A. Lamb, M. Arjovsky, O. Mastropietro, and A.-C. Courville.
336 Adversarially learned inference. In *ICLR’17*. OpenReview.net, 2017.
- 337 [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and
338 Y. Bengio. Generative adversarial nets. In *NIPS’27*, pages 2672–2680, 2014.
- 339 [13] L. Grinsztajn, E. Oyallon, and G. Varoquaux. Why do tree-based models still outperform deep
340 learning on tabular data? In *NeurIPS’22 Datasets and Benchmarks*, 2022.
- 341 [14] A. Grover and S. Ermon. Boosted generative models. In *AAAI’18*, pages 3077–3084. AAAI
342 Press, 2018.
- 343 [15] G.-E. Hinton. The forward-forward algorithm: Some preliminary investigations. *CoRR*,
344 abs/2212.13345, 2022.
- 345 [16] R.C. Holte. Very simple classification rules perform well on most commonly used datasets.
346 *MLJ*, 11:63–91, 1993.
- 347 [17] M.J. Kearns. Thoughts on hypothesis boosting, 1988. ML class project.
- 348 [18] M.J. Kearns and Y. Mansour. On the boosting ability of top-down decision tree learning
349 algorithms. In *Proc. of the 28th ACM STOC*, pages 459–468, 1996.
- 350 [19] D.-P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR’14*, 2014.
- 351 [20] D.-E. Knuth. Two notes on notation. *The American Mathematical Monthly*, 99(5):403–422,
352 1992.
- 353 [21] A. Krizhevsky, I. Sutskever, and G.-E. Hinton. ImageNet classification with deep convolutional
354 neural networks. In *NIPS’25*, pages 1106–1114, 2012.
- 355 [22] S. Lang, M. Mundt, F. Ventola, R. Peharz, and K. Kersting. Elevating perceptual sample
356 quality in pccs through differentiable sampling. In *NeurIPS 2021 Workshop on Pre-Registration*
357 in *Machine Learning, 13 December 2021, Virtual*, volume 181 of *Proceedings of Machine*
358 *Learning Research*, pages 1–25. PMLR, 2021.
- 359 [23] Q. Li and J.-S. Racine. Nonparametric estimation of distributions with categorical and continuous
360 data. *Journal of Multivariate Analysis*, 86:266–292, 2003.
- 361 [24] Y. Mansour, R. Nock, and R.-C. Williamson. Random classification noise does not defeat all
362 convex potential boosters irrespective of model choice. In *ICML’23*, 2023.

- 363 [25] P. Maréchal. On a functional operation generating convex functions, part 1: duality. *J. of
364 Optimization Theory and Applications*, 126:175–189, 2005.
- 365 [26] P. Maréchal. On a functional operation generating convex functions, part 2: algebraic properties.
366 *J. of Optimization Theory and Applications*, 126:375–366, 2005.
- 367 [27] D. McElfresh, S. Khandagale, J. Valverde, V. Prasad C, G. Ramakrishnan, M. Goldblum, and
368 C. White. When do neural nets outperform boosted trees on tabular data? In *NeurIPS’23
369 Datasets and Benchmarks*, 2023.
- 370 [28] A. Menon and C.-S. Ong. Linking losses for density ratio and class-probability estimation. In
371 *33rd ICML*, pages 304–313, 2016.
- 372 [29] R. Nock and M. Guillame-Bert. Generative trees: Adversarial and copycat. In *39th ICML*,
373 pages 16906–16951, 2022.
- 374 [30] R. Nock, A.-K. Menon, and C.-S. Ong. A scaled Bregman theorem with applications. In
375 *NIPS’29*, pages 19–27, 2016.
- 376 [31] S. Nowozin, B. Cseke, and R. Tomioka. f -GAN: training generative neural samplers using
377 variational divergence minimization. In *NIPS’29*, pages 271–279, 2016.
- 378 [32] E. Parzen. On estimation of a probability density function and mode. *The Annals of Mathematical
379 Statistics*, 33:1065–1076, 1962.
- 380 [33] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
381 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
382 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine
383 Learning Research*, 12:2825–2830, 2011.
- 384 [34] J. R. Quinlan. *C4.5 : programs for machine learning*. Morgan Kaufmann, 1993.
- 385 [35] M.-D. Reid and R.-C. Williamson. Information, divergence and risk for binary experiments.
386 *JMLR*, 12:731–817, 2011.
- 387 [36] R. Sánchez-Cauce, I. París, and F.-J. Díez. Sum-product networks: A survey. *IEEE Trans.PAMI*,
388 44(7):3821–3839, 2022.
- 389 [37] L.-J. Savage. Elicitation of personal probabilities and expectations. *J. of the Am. Stat. Assoc.*,
390 pages 783–801, 1971.
- 391 [38] M. Sugiyama and M. Kawanabe. *Machine Learning in Non-Stationary Environments - Intro-
392 duction to Covariate Shift Adaptation*. Adaptive computation and machine learning. MIT Press,
393 2012.
- 394 [39] N. Tagasovska, D. Ackerer, and T. Vatter. Copulas as high-dimensional generative models: Vine
395 copula autoencoders. In *NeurIPS’32*, pages 6525–6537, 2019.
- 396 [40] S. van Buuren and K. Groothuis-Oudshoorn. mice: Multivariate Imputation by Chained
397 Equations in R. *Journal of Statistical Software*, 45(3):1–67, 2011.
- 398 [41] A. Vergari, Y. Choi, and R. Peharz. Probabilistic circuits: Representations, inference, learning
399 and applications, 2022. *NeurIPS’22* tutorials.
- 400 [42] D.-S. Watson, K. Blesch, J. Kapar, and M.-N. Wright. Adversarial random forests for density
401 and generative modeling. In *AISTATS’23*, Proceedings of Machine Learning Research. PMLR,
402 2023.
- 403 [43] I. Witten and E. Frank. *Data Mining: Practical Machine Learning Tools and Techniques with
404 Java Implementations*. Morgan Kaufmann, 1999.
- 405 [44] C. Xiao, P. Zhong, and C. Zheng. BourGAN: Generative networks with metric embeddings. In
406 *NeurIPS’18*, pages 2275–2286, 2018.
- 407 [45] L. Xu, M. Skoulioudou, A. Cuesta-Infante, and K. Veeramachaneni. Modeling tabular data using
408 conditional GAN. In *NeurIPS’32*, pages 7333–7343, 2019.

Supplementary Material

Abstract

410 This is the Supplementary Material to Paper "Generative Forests" submitted to
 411 NeurIPS'24.

412 To differentiate with the numberings in the main file, the numbering of Theorems, etc. is letter-based
 413 (A, B, ...).

414 **Table of contents**

415 Related work	Pg 13
416	
417 Additional content	Pg 14
418	
419 Supplementary material on proofs	Pg 15
420 → Proof of Lemma 4.4	Pg 15
421 → Proof of Theorem 5.2	Pg 15
422	
423 Simpler models: ensembles of generative trees	Pg 19
424	
425 Supplementary material on experiments	Pg 21
426 → Domains	Pg 21
427 → Algorithms configuration and choice of parameters	Pg 22
428 → Interpreting our models: 'SCRUTINIZE'	Pg 24
429 → More examples of Table 1 (MF)	Pg 24
430 → The generative forest of Table 1 (MF) developed further	Pg 25
431 → Full comparisons with MICE on missing data imputation	Pg 33
432 → Experiment LIFELIKE <i>in extenso</i>	Pg 27
433 → Comparison with "the optimal generator": GEN-DISCRIM	Pg 30

434 **I Related work**

435 The typical Machine Learning (ML) problem usually contains at least three parts: (i) a training
436 algorithm minimizes (ii) a loss function to output an object whose key part is (iii) a model. The
437 main problem we are interested in is data generation as generally captured by "Generative AI". The
438 type of data we are interested in still represents one of the most prevalent form of data: tabular
439 data [6]. When it comes to tabular data, a singular phenomenon of the data world can be observed:
440 there is a fierce competition on part (iii) above, the *models*. When data has other forms, like images,
441 the ML community has generally converged to a broad idea of what the best models look like at
442 a high-level for many generic tasks: neural networks[‡]. Tabular data offers no such consensus yet,
443 even on well-defined tasks like supervised learning [13]. In fact, even on such "simple tasks" the
444 consensus is rather that there is no such consensus [27]. It is an understatement to state that in the
445 broader context of all tasks of interest to us, a sense of a truly intense competition emerges, whose
446 "gradient" clearly points towards simultaneously the most complex / expressive / tractable models, as
447 shown in [41, Slides 27, 53]. One can thus end up finding models based on trees [7, 29, 42], neural
448 networks [12, 14, 19, 45], probabilistic circuits [5, 41, 36], kernel methods [4, 32], graphical models
449 [39] (among others: note that some are in fact hybrid models).

450 So the tabular data world is blessed with a variety of possible models to solve problems like the ones
451 we are interested in, *but* – and this is another singular phenomenon of the tabular data world –, getting
452 the best solutions is not necessarily a matter of competing on size or compute. In fact, it can be the
453 opposite: striving for model simplicity or (non exclusive) lightweight tuning can substantially pay
454 off [27]. In relative terms, this phenomenon is not new in the tabular data world: it has been known
455 for more than two decades [16]. That it has endured all major revolutions in ML points to the fact
456 that lasting solutions can be conveniently addressing *all* three parts (i–iii) above at once on models,
457 algorithms and losses.

458 From this standpoint, the closest approaches to ours are [42] and [29], first and foremost because the
459 models include trees with a stochastic activation of edges to pick leaves, and a leaf-dependent data
460 generation process. While [29] learn a single tree, [42] use a way to generate data from a set of trees –
461 called an adversarial random forest – which is simple: sample a tree, and then sample an observation
462 from the tree. Hence, the distribution is a convex combination of the trees' density. This is simple but
463 it suffers several drawbacks: (i) each tree has to be accurate enough and thus big enough to model
464 "tricky" data for tree-based models (tricky can be low-dimensional, see the 2D data of Table 1, main
465 file); (ii) if leaves' samplers are simple, which is the case for [29] and our approach, it makes it tricky
466 to learn sets of simple models, such as when trees are stumps (we do not have this issue, see Table 1).
467 In our case, while our models include sets of trees, generating one observation makes use of leaves
468 in *all* trees instead of just one as in [42] (Figure 1, main file). We note that the primary goal of that
469 latter work is in fact not to generate data [42, Section 4].

470 Theoretical results that are relevant to data generation are in general scarce compared to the flurry of
471 existing methods if we omit the independent convergence rates of the toolbox often used, such as for
472 (stochastic) gradient descent. Specific convergence results are given in [42], but they are relevant to
473 statistical consistency (infinite sample) and they also rely on assumptions that are not realistic for real
474 world domains, such as Lipschitz continuity of the target density, with second derivative continuous,
475 square integrable and monotonic. The assumption made on models, that splitting probabilities on
476 trees is lowerbounded by a constant, is also impeding to model real world data.

477 In the generative trees of [29], leaf generation is the simplest possible: it is uniform. This requires big
478 trees to model real-world or tricky data. On the algorithms side, [29] introduce two training algorithms
479 in the generative adversarial networks (GAN) framework [12], thus involving the generator to train
480 but also a discriminator, which is a decision tree in [29]. The GAN framework is very convenient to
481 tackle (ii) above in the context of our tasks because it allows to tie using first principles the problem
482 of learning a density or a sampler and that of training a model (a "discriminator") to distinguish fakes
483 from real, model that parameterizes the loss optimized. An issue with neural networks is that this
484 parameterization has an uncontrollable slack unless the discriminator is extremely complex [31]. The
485 advantage of using trees as in [29] is that for such classifiers, the slack disappears because the models
486 are calibrated, so there is a tight link between training the generator and the discriminator. [29] go

[‡]Whether such a perception is caused by the models themselves or the dazzling amount of optimisation that has progressively wrapped the models, as advocated for the loss in [22], is not the focus of our paper.

Algorithm 4 ITERATIVESUPPORTUPDATE($\{\Upsilon_t\}_{t=1}^T$)

Input: Trees $\{\Upsilon_t\}_{t=1}^T$ of a GF;
Output: sampling support \mathcal{C}_s for one observation;
Step 1 : $\mathcal{C}_s \leftarrow \mathcal{X}$;
Step 2 : INIT($\{\Upsilon_t\}_{t=1}^T$);
Step 3 : **for** $t \in [T]$
 Step 2.1 : **while** Υ_t .done
 Step 2.1.1 : STARUPDATE(Υ_t , \mathcal{C}_s , R);
 return \mathcal{C}_s ;

Algorithm 5 CONCURRENTSUPPORTUPDATE

Step 1 : **while** !done
 Step 1.1 : P[accessible] // locks \mathcal{C}_s
 Step 1.2 : STARUPDATE(this, \mathcal{C}_s , R);
 Step 1.3 : V[accessible] // unlocks \mathcal{C}_s

487 further, showing that one can replace the adversarial training by a *copycat* training, involving copying
488 parts of the discriminator in the generator to speed-up training (also discussed in [15] for neural nets),
489 with strong rates on training in the original boosting model. There is however a limitation in the
490 convergence analysis of [29] as losses have to be symmetric, a property which is not desirable for
491 data generation since it ties the misclassification costs of real and fakes with no argument to do so in
492 general.

493 Our paper lifts the whole setting of [29] to models that can fit more complex densities using simpler
494 models (Table 1), using sets of trees that can be much smaller than those of [42] (Figure 1); training
495 such models is achieved by merging the two steps of copycat training into a single one where only
496 the generator is trained, furthermore keeping strong rates on training via the original boosting model,
497 all this while getting rid of the undesirable symmetry assumption of [29] for the loss at hand.

498 **II Additional content**

499 In this additional content, we provide the three ways to pick the trees in a Generative Forest to generate
500 one observation (sequential, concurrent, randomized), and then give a proof that the optimal splitting
501 threshold on a continuous variable when training a generative forest using GF BOOST is always an
502 observed value if there is one tree, but can be another value if there are more (thus highlighting some
503 technical difficulties of one wants to stick to the optimal choice of splitting).

504 **II.1 Sequentially choosing trees in a GF for the generation of observations**

505 **II.2 Concurrent generation of observations**

506 We provide a concurrent generation using Algorithm 5, which differs from Algorithm 4 (main file).
507 In concurrent generation, each tree runs concurrently algorithm UPDATESUPPORT (hence the use of
508 the Java-style `this` handler), with an additional global variables (in addition to \mathcal{C}_s , initialized to \mathcal{X}):
509 a Boolean semaphore `accessible` implementing a lock, whose value 1 means \mathcal{C} is available for an
510 update (and otherwise it is locked by a tree in Steps 1.2/1.3 in the set of trees of the GF). We assume
511 that INIT has been run beforehand (eventually locally).

512 **II.3 Randomized generation of observations**

513 We provide a general randomized choice of the sequence of trees for generation, in Algorithm 6.

Algorithm 6 RANDOMIZEDSUPPORTUPDATE

Step 1 : INIT($\{\Upsilon_t\}_{t=1}^T$); $\mathbb{I} \leftarrow \{1, 2, \dots, T\}$;
Step 2 : **while** ($\mathbb{I} \neq \emptyset$)
 Step 2.1 : $i \leftarrow \text{RANDOM}(\mathbb{I})$;
 Step 2.2 : STARUPDATE($\Upsilon_i, \mathcal{C}_s, \mathbf{R}$);
 Step 2.3 : **if** ($\Upsilon_i.\text{done}$) **then** $\mathbb{I} \leftarrow \mathbb{I} \setminus \{i\}$;

514 **III Supplementary material on proofs**

515 **III.1 Proof of Lemma 4.4**

516 Given sequence \mathbf{v} of dimension $\dim(\mathbf{v})$, denote $\{\mathcal{C}_j\}_{j \in [1 + \dim(\mathbf{v})]}$ the sequence of subsets of the
517 domain appearing in the parameters of UPDATE SUPPORT through sequence \mathbf{v} , to which we add a last
518 element, \mathcal{C}_s (and its first element is \mathcal{X}). If we let \mathcal{X}_j denote the support of the corresponding star node
519 whose Bernoulli event is triggered at index j in sequence \mathbf{v} (for example, $\mathcal{X}_1 = \mathcal{X}$), then we easily
520 get

$$\mathcal{C}_j \subseteq \mathcal{X}_j, \forall j \in [\dim(\mathbf{v})], \quad (5)$$

521 indicating the Bernoulli probability in Step 1.1 of STARUPDATE is always defined. We then compute

$$p_{\mathbf{G}}[\mathcal{C}_s | \mathbf{v}] = p_{\mathbf{G}}(\cap_j \mathcal{C}_j) \quad (6)$$

$$= \prod_{j=1}^N p_{\mathbf{G}}[\mathcal{C}_{j+1} | \mathcal{C}_j] \quad (7)$$

$$= \prod_{j=1}^N \frac{p_{\mathbf{R}}[\mathcal{C}_{i+1}]}{p_{\mathbf{R}}[\mathcal{C}_i]} \quad (8)$$

$$= \frac{p_{\mathbf{R}}[\mathcal{C}_{N+1}]}{p_{\mathbf{R}}[\mathcal{C}_0]} \quad (9)$$

$$= \frac{p_{\mathbf{R}}[\mathcal{C}_s]}{p_{\mathbf{R}}[\mathcal{X}]} \quad (10)$$

522 (7) holds because updating the support generation is Markovian in a GF, (8) holds because of Step
523 3.2 and (9) is a simplification through cancelling terms in products.

524 **III.2 Proof of Theorem 5.2**

525 Notations: we iteratively learn generators $\mathbf{G}_0, \mathbf{G}_1, \dots, \mathbf{G}_J$ where \mathbf{G}_0 is just a root. We let $\mathcal{P}(\mathbf{G}_j)$
526 denote the partition induced by the set of trees of \mathbf{G}_j , recalling that each element is the (non-empty)
527 intersection of the support of a set of leaves, one for each tree (for example, $\mathcal{P}(\mathbf{G}_0) = \{\mathcal{X}\}$). The
528 criterion minimized to build \mathbf{G}_{j+1} from \mathbf{G}_j is

$$\underline{\mathbb{L}}(\mathbf{G}_j) \doteq \sum_{\mathcal{C} \in \mathcal{P}(\mathbf{G}_j)} p_{\mathbf{M}}[\mathcal{C}] \cdot \underline{L}\left(\frac{\pi p_{\mathbf{R}}[\mathcal{C}]}{p_{\mathbf{M}}[\mathcal{C}]}\right). \quad (11)$$

529 The proof entails fundamental notions on loss functions, models and boosting. We start with loss
530 functions. A key function we use is

$$g^\pi(t) \doteq (\pi t + 1 - \pi) \cdot \underline{L}\left(\frac{\pi t}{\pi t + 1 - \pi}\right), \forall t \in \mathbb{R}_+,$$

531 which is concave [35, Appendix A.3].

Lemma A.

$$\mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_{J-1}) - \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_J) = \underline{\mathbb{L}}(\mathbf{G}_{J-1}) - \underline{\mathbb{L}}(\mathbf{G}_J). \quad (12)$$

532 *Proof.* We make use of the following important fact about GFs:

533 **(F1)** For any $\mathcal{X}' \in \mathcal{P}_J$, $\int_{\mathcal{X}'} d\mathbf{G}_J = \int_{\mathcal{X}'} d\mathbf{R}$.

534 We have from the proof of [29, Lemma 5.6],

$$\mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_J) = \sum_{\mathcal{X}' \in \mathcal{P}_J} \left\{ g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}'} \left[\frac{d\mathbf{G}_J}{d\mathbf{U}} \right] \right) - \mathbb{E}_{\mathbf{U}|\mathcal{X}'} \left[g^\pi \left(\frac{d\mathbf{R}}{d\mathbf{U}} \right) \right] \right\}. \quad (13)$$

535 We note that **(F1)** implies (13) is just a sum of slacks of Jensen's inequality. We observe

$$\begin{aligned} & \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_{J-1}) - \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_J) \\ &= \sum_{\mathcal{X}' \in \mathcal{P}_{J-1}} \left\{ g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}'} \left[\frac{d\mathbf{G}_{J-1}}{d\mathbf{U}} \right] \right) - \mathbb{E}_{\mathbf{U}|\mathcal{X}'} \left[g^\pi \left(\frac{d\mathbf{R}}{d\mathbf{U}} \right) \right] \right\} \\ &\quad - \sum_{\mathcal{X}' \in \mathcal{P}_J} \left\{ g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}'} \left[\frac{d\mathbf{G}_J}{d\mathbf{U}} \right] \right) - \mathbb{E}_{\mathbf{U}|\mathcal{X}'} \left[g^\pi \left(\frac{d\mathbf{R}}{d\mathbf{U}} \right) \right] \right\} \\ &= \sum_{(\mathcal{X}_t, \mathcal{X}_f) \in \mathcal{P}_J^{\text{split}}} \left\{ g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}_t \cup \mathcal{X}_f} \left[\frac{d\mathbf{G}_{J-1}}{d\mathbf{U}} \right] \right) - \sum_{v \in \{t, f\}} g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}_v} \left[\frac{d\mathbf{G}_J}{d\mathbf{U}} \right] \right) \right\}, \end{aligned}$$

536 where $\mathcal{P}_J^{\text{split}}$ contains all couples $(\mathcal{X}_t, \mathcal{X}_f)$ such that $\mathcal{X}_t, \mathcal{X}_f \in \mathcal{P}_J$ and $\mathcal{X}_t \cup \mathcal{X}_f \in \mathcal{P}_{J-1}$. These unions
537 were subsets of the partition of \mathcal{X} induced by the set of trees of \mathbf{G}_{J-1} and that were cut by the
538 predicate p put at the leaf λ_* that created \mathbf{G}_J from \mathbf{G}_{J-1} . To save space, denote $\mathcal{X}_{tf} \doteq \mathcal{X}_t \cup \mathcal{X}_f$.

$$\begin{aligned} & g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}_{tf}} \left[\frac{d\mathbf{G}_{J-1}}{d\mathbf{U}} \right] \right) - \sum_{v \in \{t, f\}} g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}_v} \left[\frac{d\mathbf{G}_J}{d\mathbf{U}} \right] \right) \\ &= \left(\pi \int_{\mathcal{X}_{tf}} d\mathbf{G}_{J-1} + (1 - \pi) \int_{\mathcal{X}_{tf}} d\mathbf{U} \right) \cdot \underline{L} \left(\frac{\pi \int_{\mathcal{X}_{tf}} d\mathbf{G}_{J-1}}{\pi \int_{\mathcal{X}_{tf}} d\mathbf{G}_{J-1} + (1 - \pi) \int_{\mathcal{X}_{tf}} d\mathbf{U}} \right) \\ &\quad - \left(\pi \int_{\mathcal{X}_t} d\mathbf{G}_J + (1 - \pi) \int_{\mathcal{X}_t} d\mathbf{U} \right) \cdot \underline{L} \left(\frac{\pi \int_{\mathcal{X}_t} d\mathbf{G}_J}{\pi \int_{\mathcal{X}_t} d\mathbf{G}_J + (1 - \pi) \int_{\mathcal{X}_t} d\mathbf{U}} \right) \\ &\quad - \left(\pi \int_{\mathcal{X}_f} d\mathbf{G}_J + (1 - \pi) \int_{\mathcal{X}_f} d\mathbf{U} \right) \cdot \underline{L} \left(\frac{\pi \int_{\mathcal{X}_f} d\mathbf{G}_J}{\pi \int_{\mathcal{X}_f} d\mathbf{G}_J + (1 - \pi) \int_{\mathcal{X}_f} d\mathbf{U}} \right). \quad (14) \end{aligned}$$

539 We now work on (14). Using **(F1)**, we note $\int_{\mathcal{X}_{tf}} d\mathbf{G}_{J-1} = \int_{\mathcal{X}_{tf}} d\mathbf{R}$ since $\mathcal{X}_{tf} \in \mathcal{P}_{J-1}$. Similarly,
540 $\int_{\mathcal{X}_v} d\mathbf{G}_J = \int_{\mathcal{X}_v} d\mathbf{R}$, $\forall v \in \{t, f\}$, so we make appear \mathbf{R} in (14) and get:

$$\begin{aligned} & g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}_{tf}} \left[\frac{d\mathbf{G}_{J-1}}{d\mathbf{U}} \right] \right) - \sum_{v \in \{t, f\}} g^\pi \left(\mathbb{E}_{\mathbf{U}|\mathcal{X}_v} \left[\frac{d\mathbf{G}_J}{d\mathbf{U}} \right] \right) \\ &= p(tf) \cdot \left\{ \underline{L} \left(\frac{\pi \int_{\mathcal{X}_{tf}} d\mathbf{R}}{p(tf)} \right) - \frac{p(t)}{p(tf)} \cdot \underline{L} \left(\frac{\pi \int_{\mathcal{X}_t} d\mathbf{R}}{p(t)} \right) - \frac{p(f)}{p(tf)} \cdot \underline{L} \left(\frac{\pi \int_{\mathcal{X}_f} d\mathbf{R}}{p(f)} \right) \right\}, \quad (15) \end{aligned}$$

541 where we let for short

$$p(v) \doteq \pi \int_{\mathcal{X}_v} d\mathbf{R} + (1 - \pi) \int_{\mathcal{X}_v} d\mathbf{U}, \forall v \in \{t, v, tv\}. \quad (16)$$

542 We finally get

$$\begin{aligned} & \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_{J-1}) - \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_J) \\ &= \sum_{\mathcal{C} \in \mathcal{P}_{J-1}} p_M[\mathcal{C}] \cdot \underline{L} \left(\frac{\pi p_R[\mathcal{C}]}{p_M[\mathcal{C}]} \right) - \sum_{\mathcal{C} \in \mathcal{P}_{J-1}} p_M[\mathcal{C}] \cdot \sum_{v \in \{t, f\}} \frac{p_M[\mathcal{C}_{p_{J-1}, v}]}{p_M[\mathcal{C}]} \cdot \underline{L} \left(\frac{\pi p_R[\mathcal{C}_{p_{J-1}, v}]}{p_M[\mathcal{C}_{p_{J-1}, v}]} \right) \\ &= \underline{L}(\mathbf{G}_{J-1}) - \underline{L}(\mathbf{G}_J), \end{aligned}$$

543 as claimed. The last identity comes from the fact that the contribution to $\underline{L}(.)$ is the same outside
544 \mathcal{P}_{J-1} for both \mathbf{G}_{J-1} and \mathbf{G}_J . \square

Algorithm 7 MODABOOST($\mathbf{R}, \ell, \text{WL}, J, T$)

Input: measure \mathbf{R} , SPD loss ℓ , weak learner WL , iteration number $J \geq 1$, number of trees T ;
Output: PLM (partition-linear model) H_J with T trees;

Step 1 : $\Gamma_0 \doteq \{\Upsilon_t\}_{t=1}^T \leftarrow \{\text{root}, \text{root}, \dots, \text{root}\}; \mathcal{P}(\Gamma_0) \leftarrow \{\mathcal{X}\};$

Step 2 : **for** $j = 1, 2, \dots, J$

- Step 2.1 : pick $t \in [T];$
- Step 2.2 : pick $\lambda \in \Lambda(\Upsilon_t);$
- Step 2.3 : $h_j \leftarrow \text{WL}(\mathcal{X}_\lambda);$
 $\quad // \text{call to the weak learner for some } h_j \text{ of the type } h_j = 1_{p(.)} \cdot K_j,$
 $\quad // K_j \text{ constant, } p(.) \text{ a splitting predicate (e.g. } x_i \geq a)$
- Step 2.4 : **for** $\mathcal{C} \in \mathcal{P}_\lambda(\Gamma_{j-1}),$

$$\mathcal{C}^p \leftarrow \mathcal{C} \cap \{x : p(x) \text{ is true}\}, \quad (17)$$

$$\alpha(\mathcal{C}^p) \leftarrow (-\underline{L}') \left(\frac{\pi p_{\mathbf{R}}[\mathcal{C}^p]}{p_{\mathbf{M}}[\mathcal{C}^p]} \right), \quad (18)$$

$$\mathcal{C}^{-p} \leftarrow \mathcal{C} \cap \{x : p(x) \text{ is false}\}, \quad (19)$$

$$\alpha(\mathcal{C}^{-p}) \leftarrow -(-\underline{L}') \left(\frac{\pi p_{\mathbf{R}}[\mathcal{C}^{-p}]}{p_{\mathbf{M}}[\mathcal{C}^{-p}]} \right). \quad (20)$$

Step 2.5 : // update $\mathcal{P}(\Gamma_j):$

$$\mathcal{P}(\Gamma_j) \leftarrow (\mathcal{P}(\Gamma_{j-1}) \setminus \mathcal{P}_\lambda(\Gamma_{j-1})) \cup (\cup_{\mathcal{C} \in \mathcal{P}_\lambda(\Gamma_{j-1}): \mathcal{C}^p \neq \emptyset} \mathcal{C}^p) \cup (\cup_{\mathcal{C} \in \mathcal{P}_\lambda(\Gamma_{j-1}): \mathcal{C}^{-p} \neq \emptyset} \mathcal{C}^{-p})$$

Step 2.6 : split leaf λ at Υ_t using predicate p ;

return $H_J(\mathbf{x}) \doteq \alpha(\mathcal{C}(\mathbf{x}))$ with $\mathcal{C}(\mathbf{x}) \doteq \mathcal{C} \in \mathcal{P}(\Gamma_J)$ such that $\mathbf{x} \in \mathcal{C};$

545 We now come to models and boosting. Suppose we have t trees $\{\Upsilon_1, \Upsilon_2, \dots, \Upsilon_t\}$ in \mathbf{G}_j . We want
546 to split leaf $\lambda \in \Upsilon_t$ into two new leaves, λ_f, λ_t . Let $\mathcal{P}_\lambda(\mathbf{G}_j)$ denote the subset of $\mathcal{P}(\mathbf{G}_j)$ containing
547 only the subsets defined by intersection with the support of λ , \mathcal{X}_λ . The criterion we minimize can be
548 reformulated as

$$\underline{L}(\mathbf{G}_j) = \sum_{\lambda \in \Lambda(\Upsilon)} p_{\mathbf{M}}[\mathcal{X}_\lambda] \cdot \sum_{\mathcal{C} \in \mathcal{P}_\lambda(\mathbf{G}_j)} \frac{p_{\mathbf{M}}[\mathcal{C}]}{p_{\mathbf{M}}[\mathcal{X}_\lambda]} \cdot \underline{L} \left(\frac{\pi p_{\mathbf{R}}[\mathcal{C}]}{p_{\mathbf{M}}[\mathcal{C}]} \right), \Upsilon \in \{\Upsilon_1, \Upsilon_2, \dots, \Upsilon_t\}.$$

549 Here, Υ is any tree in the set of trees, since $\cup_{\lambda \in \Lambda(\Upsilon)} \mathcal{P}_\lambda(\mathbf{G}_j)$ covers the complete partition of \mathcal{X}
550 induced by $\{\Upsilon_1, \Upsilon_2, \dots, \Upsilon_t\}$. If we had a single tree, the inner sum would disappear since we would
551 have $\mathcal{P}_\lambda(\mathbf{G}_j) = \{\mathcal{X}_\lambda\}$ and so one iteration would split one of these subsets. In our case however,
552 with a set of trees, we still split \mathcal{X}_λ but the impact on the reduction of \underline{L} can be substantially better
553 as we may simultaneously split as many subsets as there are in $\mathcal{P}_\lambda(\mathbf{G}_j)$. The reduction in \underline{L} can be
554 obtained by summing all reductions to which contribute each of the subsets.

555 To analyze it, we make use of a reduction to the MODABOOST algorithm of [24]. We present the
556 algorithm in Algorithm 7. The original MODABOOST algorithm trains *partition-linear models*, i.e.
557 models whose output is defined by a sum of reals over a set of partitions to which the observation to be
558 classified belongs to. Training means then both learning the organisation of partitions and the reals –
559 which are just the output of a weak classifier given by a weak learner, times a leveraging constant com-
560 puted by MODABOOST. As in the classical boosting literature, the original MODABOOST includes
561 the computation and updates of *weights*.

562 In our case, the structure of partition learned is that of a set of trees, each weak classifier h is of the
563 form

$$h(\mathbf{x}) \doteq K \cdot p(\mathbf{x}),$$

564 where K is a real and p is a Boolean predicate, usually doing an axis-parallel split of \mathcal{X} on an observed
565 feature, such as $p(\mathbf{x}) \equiv 1_{x_i \geq a}$.

566 From now on, it will be interesting to dissociate our generator \mathbf{G}_j – which includes a model structure
567 and an algorithm to generate data from this structure – from its set of trees – *i.e.* its model structure –
568 which we denote $\Gamma_j \doteq \{\Upsilon_t\}_{t=1}^T$. MODABOOST learns both Γ_j but also predictions for each possible

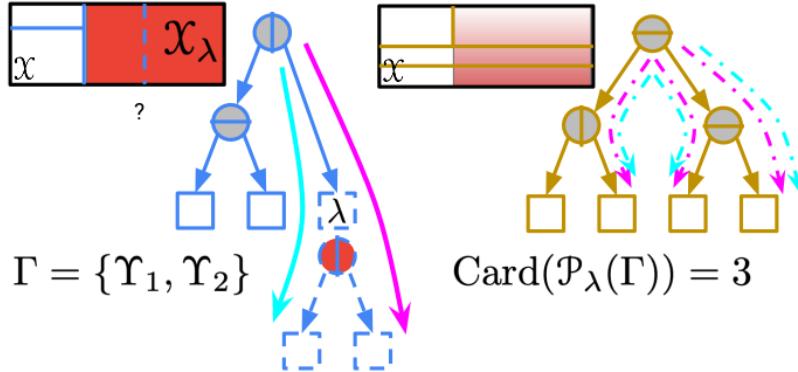


Figure III.1: Learning a set of two trees Γ with MODABOOST. If we want a split at leaf λ indicated, then it would take two steps (but a *single* application of the weak learning assumption, [24, Lemma 5]) of the original MODABOOST to learn Υ_1 alone [24, Section 4]. Each step, figured with a plain arrow, ends with adding a new leaf. In our case however, we have to take into account the partition induced by the leaves of Υ_2 , which results in considering not one but three subsets in \mathcal{P}_λ , thus adding not one but three weak hypotheses simultaneously (one for each of the dashed arrows on Υ_2). Thanks to Jensen's inequality, a guaranteed decrease of the loss at hand can be obtained by a single application of the weak learning assumption at λ , just like in the original MODABOOST.

569 outcome, predictions that we shall not use since the loss we are minimizing can be made to depend
 570 only on the model structure Γ_j . Given the simplicity of the weak classifiers h and the particular nature
 571 of the partitions learned, the original MODABOOST of [24] can be simplified to our presentation
 572 in Algorithm 7. Among others, the simplification discards the weights from the presentation of the
 573 algorithm. What we show entangles two objectives on models and loss functions as we show that

574 MODABOOST learns the same structure Γ as our GF.BOOST, and yields a guaranteed decrease of
 575 $\underline{\mathbb{L}}(\mathbf{G})$,

576 and it is achieved via the following Theorem.

577 **Theorem B.** MODABOOST greedily minimizes the following loss function:

$$\underline{\mathbb{L}}(\Gamma_j) = \sum_{\lambda \in \Lambda(\Upsilon)} p_{\text{M}}[\mathcal{X}_\lambda] \cdot \sum_{\mathcal{C} \in \mathcal{P}_\lambda(\Gamma_j)} \frac{p_{\text{M}}[\mathcal{C}]}{p_{\text{M}}[\mathcal{X}_\lambda]} \cdot \underline{\mathbb{L}}\left(\frac{\pi p_{\text{R}}[\mathcal{C}]}{p_{\text{M}}[\mathcal{C}]}\right), \Upsilon \in \Gamma_j.$$

578 Furthermore, suppose there exists $\gamma > 0, \kappa > 0$ such that at each iteration j , the predicate p splits
 579 \mathcal{X}_λ of leaf $\lambda \in \Lambda(\Upsilon)$ into \mathcal{X}_λ^p and \mathcal{X}_λ^{-p} (same nomenclature as in (17), (19)) such that:

$$p_{\text{M}}[\mathcal{X}_\lambda] \geq \frac{1}{\text{Card}(\Lambda(\Upsilon))}, \quad (22)$$

$$\left| \frac{\int_{\mathcal{X}_\lambda^p} d\mathbf{R}}{\int_{\mathcal{X}_\lambda} d\mathbf{R}} - \frac{\int_{\mathcal{X}_\lambda^{-p}} d\mathbf{U}}{\int_{\mathcal{X}_\lambda} d\mathbf{U}} \right| \geq \gamma, \quad (23)$$

$$\min \left\{ \frac{\pi p_{\text{R}}[\mathcal{X}_\lambda]}{p_{\text{M}}[\mathcal{X}_\lambda]}, 1 - \frac{\pi p_{\text{R}}[\mathcal{X}_\lambda]}{p_{\text{M}}[\mathcal{X}_\lambda]} \right\} \geq \kappa. \quad (24)$$

580 Then we have the following guaranteed slack between two successive models:

$$\underline{\mathbb{L}}(\Gamma_{j+1}) - \underline{\mathbb{L}}(\Gamma_j) \leq -\frac{\kappa \gamma^2 \kappa^2}{8 \text{Card}(\Lambda(\Upsilon))}, \quad (25)$$

581 where κ is such that $0 < \kappa \leq \inf\{\ell'_{-1} - \ell'_1\}$.

582 **Proof sketch** The loss function comes directly from [24, Lemma 7]. At each iteration,
 583 MODABOOST makes a number of updates that guarantee, once Step 2.4 is completed (because

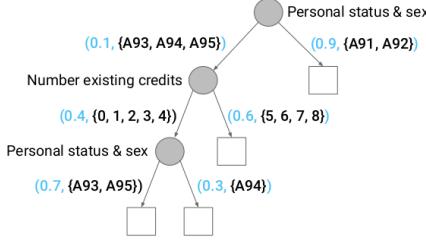


Figure IV.1: A generative tree (GT) associated to UCI German Credit.

584 the elements in $\mathcal{P}_\lambda(\Gamma_j)$ are disjoint

$$\underline{\mathbb{L}}(\Gamma_{j+1}) - \underline{\mathbb{L}}(\Gamma_j) \leq -\frac{\kappa}{2} \cdot \sum_{\mathcal{C} \in \mathcal{P}_\lambda(\Gamma_j)} p_{\textcolor{red}{M}}[\mathcal{C}] \cdot \mathbb{E}_{\textcolor{red}{M}|\mathcal{C}} [(w_{j+1} - w_j)^2] \quad (26)$$

$$= -\frac{\kappa}{2} \cdot p_{\textcolor{red}{M}}[\mathcal{X}_\lambda] \mathbb{E}_{\textcolor{red}{M}|\mathcal{X}_\lambda} [(w_{j+1} - w_j)^2], \quad (27)$$

585 where the weights are given in [24]. Each expression in the summand of (26) is exactly the guarantee
 586 of [24, ineq. before (69)]; all such expressions are not important; what is more important is (26): all
 587 steps occurring in Step 2.4 are equivalent to a single step of MODABOOST carried out over the whole
 588 \mathcal{X}_λ . Overall, the number of "aggregated" steps match the counter j in MODABOOST. We then just
 589 have to reuse the proof of [24, Theorem B], which implies, in lieu of their [24, eq. (74)]

$$\underline{\mathbb{L}}(\Gamma_{j+1}) - \underline{\mathbb{L}}(\Gamma_j) \leq -2\kappa \cdot p_{\textcolor{red}{M}}[\mathcal{X}_\lambda] \cdot \left(\frac{1}{2} \cdot \left(\frac{\int_{\mathcal{X}_\lambda^p} d\mathbf{R}}{\int_{\mathcal{X}_\lambda} d\mathbf{R}} - \frac{\int_{\mathcal{X}_\lambda^p} d\mathbf{U}}{\int_{\mathcal{X}_\lambda} d\mathbf{U}} \right) \right)^2 \cdot \underline{L}^{\text{sq}} \left(\frac{\pi p_{\textcolor{red}{R}}[\mathcal{X}_\lambda]}{p_{\textcolor{red}{M}}[\mathcal{X}_\lambda]} \right)^2 \quad (28)$$

590 with $\underline{L}^{\text{sq}}(u) \doteq u(1-u)$. Noting $2\underline{L}^{\text{sq}}(u) \geq \min\{u, 1-u\}$, we then use (22) – (24), which yields
 591 the statement of the Theorem. ■

592 As a consequence, if we assume that the total number of boosting iterations J is a multiple of the
 593 number of trees T , it comes from [24, eq. (29)] that after J iterations, we have

$$\underline{\mathbb{L}}(\Gamma_J) - \underline{\mathbb{L}}(\Gamma_0) \leq -\frac{\kappa\gamma^2\kappa^2}{8} \cdot T \log \left(1 + \frac{J}{T} \right). \quad (29)$$

594 Using Lemma A and the fact that the induction of the sets of trees in MODABOOST is done in the
 595 same way as the induction of the set of trees of our generator \mathbf{G} in GF.BOOST, we get:

$$\begin{aligned} \pi \cdot \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_J) &= \pi \cdot \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_0) + (\underline{\mathbb{L}}(\mathbf{G}_J) - \underline{\mathbb{L}}(\mathbf{G}_0)) \\ &= \pi \cdot \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_0) + (\underline{\mathbb{L}}(\Gamma_J) - \underline{\mathbb{L}}(\Gamma_0)) \\ &\leq \pi \cdot \mathbb{D}_\ell(\mathbf{R}, \mathbf{G}_0) - \frac{\kappa\gamma^2\kappa^2}{8} \cdot T \log \left(1 + \frac{J}{T} \right), \end{aligned}$$

596 as claimed.

597 IV Simpler models: ensembles of generative trees

598 Storing a GF requires keeping information about the empirical measure \mathbf{R} to compute the branching
 599 probability in Step 1 of STARUPDATE. This does not require to store the full training sample, but
 600 requires at least an index table recording the $\{\text{leaves}\} \times \{\text{observations}\}$ association, for a storage cost
 601 in between $\Omega(m)$ and $O(mT)$ where m is the size of the training sample. There is a simple way to
 602 get rid of this constraint and approximate the GF by a set of *generative trees* (GTs) of [29].

603 **Models** Simply put, a GT is a simplified equivalent representation of a generative forest with 1 tree
 604 only. In this case, the branching probabilities in Step 1 of STARUPDATE depend only on the tree's
 605 star node position. Thus, instead of recomputing them from scratch each time an observation needs to
 606 be generated, we compute them beforehand and use them to label the arcs in addition to the features'
 607 domains, as shown in Figure IV.1. This representation is equivalent to that of generative trees [29]. If

608 we have several trees, we do this process independently for each tree and end up with an *ensemble of*
 609 *generative trees* (EOGT). The additional memory footprint for storing probabilities ($O(\sum_i |\Lambda(\Upsilon_i)|)$)
 610 is offset by the fact that we do not have anymore to store associations between leaves and observations
 611 for generative forests, for a potentially significant reduction in storing size. However, we cannot use
 612 anymore STARUPDATE as is for data generation since we cannot compute exactly the probabilities in
 613 Step 1. Two questions need to be addressed: is it possible to use an ensemble of generative trees to
 614 generate data with good approximation guarantees (and if so, how) and of course how do we train
 615 such models.

616 **Data generation** We propose a simple solution based on how well an EOGT can approximate a
 617 generative forest. It relies on a simple assumption about \mathbf{R} and \mathcal{X} . Taking as references the parameters
 618 in STARUPDATE, we now make two simplifications on notations, first replacing $\Upsilon.\nu_v^*$ by ν_v^* (tree
 619 implicit), and then using notation $\mathcal{C}_v \doteq \mathcal{X}_{\nu_v^*} \cap \mathcal{C}$ for any $v \in \{t, f\}$, the reference to the tree / star
 620 node being implicit. The assumption is as follows.

621 **Assumption A.** *There exists $\varkappa \in (0, \infty)$ such that at any Step 1 of STARUPDATE, for any $v \in \{t, f\}$,
 622 we have*

$$\frac{p_{\mathbf{R}}[\mathcal{C}_v | \mathcal{X}_{\nu_v^*}]}{p_{\mathbf{U}}[\mathcal{C}_v | \mathcal{X}_{\nu_v^*}]} \in [\exp(-\varkappa), \exp(\varkappa)]. \quad (30)$$

623 A simple condition to ensure the existence of \varkappa is a requirement weaker than (c) in Definition 5.1:
 624 at all Steps 1 of STARUPDATE, $p_{\mathbf{R}}[\mathcal{C}_t | \mathcal{C}] \in (0, 1)$, which also guarantees $p_{\mathbf{R}}[\mathcal{C}_f | \mathcal{C}] \in (0, 1)$ and
 625 thus postulates that the branching in Step 1 of STARUPDATE never reduces to one choice only. The
 626 next Lemma shows that it is indeed possible to combine the generation of an ensemble of generative
 627 trees with good approximation properties, and provides a simple algorithm to do so, which reduces
 628 to running STARUPDATE with a specific approximation to branching probabilities in Step 1. For
 629 any GF \mathbf{G} and $\mathcal{C} \in \mathcal{P}(\mathbf{G})$, $\text{depth}_{\mathbf{G}}(\mathcal{C})$ is the sum of depths of the leaves in each tree whose support
 630 intersection is \mathcal{C} . We need the definition of the *expected depth* of \mathbf{G} .

631 **Definition B.** *The expected depth of GF \mathbf{G} is $\overline{\text{depth}}(\mathbf{G}) \doteq \sum_{\mathcal{C} \in \mathcal{P}(\mathbf{G})} p_{\mathbf{G}}[\mathcal{C}] \cdot \text{depth}_{\mathbf{G}}(\mathcal{C})$.*

632 $\overline{\text{depth}}(\mathbf{G})$ represents the expected complexity to sample an observation (see also Lemma C).

633 **Lemma C.** *In Step 1 of STARUPDATE, suppose $p_{\mathbf{R}}[\mathcal{C}_t \cap \mathcal{C} | \mathcal{C}]$ is replaced by*

$$\hat{p}_{\nu^*} \doteq \frac{p_{\mathbf{U}}[\mathcal{C}_t | \mathcal{X}_{\nu_t^*}] \cdot p_{\nu^*}}{p_{\mathbf{U}}[\mathcal{C}_t | \mathcal{X}_{\nu_t^*}] \cdot p_{\nu^*} + p_{\mathbf{U}}[\mathcal{C}_f | \mathcal{X}_{\nu_f^*}] \cdot (1 - p_{\nu^*})}, \quad (31)$$

634 with $p_{\nu^*} \doteq p_{\mathbf{R}}[\mathcal{X}_{\nu_t^*} | \mathcal{X}_{\nu^*}]$ the probability labeling arc (ν^*, ν_t) in its generative tree. Under Assump-
 635 tion A, if we denote \mathbf{G} the initial GF and $\hat{\mathbf{G}}$ the EOGT.P using (31), the following bound holds on the
 636 KL divergence between \mathbf{G} and $\hat{\mathbf{G}}$:

$$\text{KL}(\mathbf{G} \| \hat{\mathbf{G}}) \leq 2\varkappa \cdot \overline{\text{depth}}(\mathbf{G}), \quad \text{where } \overline{\text{depth}}(\mathbf{G}) \text{ is given in Definition B.} \quad (32)$$

637 *Proof.* Because $\mathcal{C} \subseteq \mathcal{X}_{\Upsilon.\nu^*}$ at any call of STARUPDATE, we have

$$\begin{aligned} & p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_t^*} \cap \mathcal{C} | \mathcal{C}] \\ &= p_{\mathbf{R}}[\mathcal{C}_t | \mathcal{C}] \\ &= \frac{p_{\mathbf{R}}[\mathcal{C}_t]}{p_{\mathbf{R}}[\mathcal{C}_t] + p_{\mathbf{R}}[\mathcal{C}_f]} \end{aligned} \quad (33)$$

$$\leq \frac{\exp(\varkappa)p_{\mathbf{U}}[\mathcal{C}_t | \mathcal{X}_{\Upsilon.\nu_t^*}]p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_t^*}]}{\exp(-\varkappa)p_{\mathbf{U}}[\mathcal{C}_t | \mathcal{X}_{\Upsilon.\nu_t^*}]p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_t^*}] + \exp(-\varkappa)p_{\mathbf{U}}[\mathcal{C}_f | \mathcal{X}_{\Upsilon.\nu_f^*}]p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_f^*}]} \quad (34)$$

$$= \exp(2\varkappa) \cdot \frac{p_{\mathbf{U}}[\mathcal{C}_t | \mathcal{X}_{\Upsilon.\nu_t^*}]p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_t^*} | \mathcal{X}_{\Upsilon.\nu^*}]}{p_{\mathbf{U}}[\mathcal{C}_t | \mathcal{X}_{\Upsilon.\nu_t^*}]p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_t^*} | \mathcal{X}_{\Upsilon.\nu^*}] + p_{\mathbf{U}}[\mathcal{C}_f | \mathcal{X}_{\Upsilon.\nu_f^*}]p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_f^*} | \mathcal{X}_{\Upsilon.\nu^*}]} \quad (35)$$

638 (34) is due to the fact that, since $\mathcal{C} \subseteq \mathcal{X}_{\Upsilon.\nu^*}$, $p_{\mathbf{R}}[\mathcal{C}_v] = p_{\mathbf{R}}[\mathcal{C}_v \cap \mathcal{X}_{\Upsilon.\nu_v^*}] = p_{\mathbf{R}}[\mathcal{C}_v | \mathcal{X}_{\Upsilon.\nu_v^*}]p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_v^*}]$,
 639 and then using (30). (35) is obtained by dividing numerator and denominator by $p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_t^*}] +$
 640 $p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu_f^*}] = p_{\mathbf{R}}[\mathcal{X}_{\Upsilon.\nu^*}]$.

641 \mathbf{G} and $\hat{\mathbf{G}}$ satisfy $\mathcal{P}(\mathbf{G}) = \mathcal{P}(\hat{\mathbf{G}})$ so

$$\text{KL}(\mathbf{G}\|\hat{\mathbf{G}}) = \int d\mathbf{G} \log \frac{d\mathbf{G}}{d\hat{\mathbf{G}}} \quad (36)$$

$$= \sum_{\mathcal{C} \in \mathcal{P}(\mathbf{G})} p_{\mathbf{G}}[\mathcal{C}] \log \frac{p_{\mathbf{G}}[\mathcal{C}]}{p_{\hat{\mathbf{G}}}[\mathcal{C}]} \quad (37)$$

642 Finally, $p_{\mathbf{G}}[\mathcal{C}]$ and $p_{\hat{\mathbf{G}}}[\mathcal{C}]$ are just the product of the branching probabilities in any admissible
643 sequence. If we use the same admissible sequence in both generators, we can write $p_{\mathbf{G}}[\mathcal{C}] = \prod_{j=1}^{n(\mathcal{C})} p_j$
644 and $p_{\hat{\mathbf{G}}}[\mathcal{C}] = \prod_{j=1}^{n(\mathcal{C})} \hat{p}_j$, and (35) directly yields $p_j \leq \exp(2\varkappa) \cdot \hat{p}_j, \forall j \in [n(\mathcal{C})]$, $n(\mathcal{C})$ being a
645 shorthand for $\text{depth}_{\mathbf{G}}(\mathcal{C}) = \text{depth}_{\hat{\mathbf{G}}}(\mathcal{C})$ (see main file). So, for any $\mathcal{C} \in \mathcal{P}(\mathbf{G})$,

$$\frac{p_{\mathbf{G}}[\mathcal{C}]}{p_{\hat{\mathbf{G}}}[\mathcal{C}]} \leq \exp(2\varkappa \cdot \text{depth}_{\mathbf{G}}(\mathcal{C})), \quad (38)$$

646 and finally, replacing back $n(\mathcal{C})$ by notation $\text{depth}(\mathcal{C})$,

$$\begin{aligned} \text{KL}(\mathbf{G}\|\hat{\mathbf{G}}) &\leq 2\varkappa \cdot \sum_{\mathcal{C} \in \mathcal{P}(\mathbf{G})} p_{\mathbf{G}}[\mathcal{C}] \cdot \text{depth}(\mathcal{C}) \\ &= 2\varkappa \cdot \overline{\text{depth}}(\mathbf{G}), \end{aligned} \quad (39)$$

647 as claimed. \square

648 Note the additional leverage for training and generation that stems from Assumption A, not just in
649 terms of space: computing $p_{\mathbf{U}}[\mathcal{C}_v | \mathcal{X}_{\nu_v^*}]$ is $O(d)$ and does not necessitate data so the computation of
650 key conditional probabilities (31) drops from $\Omega(m)$ to $O(d)$ for training and generation in an EOGT.
651 Lemma C provides the change in STARUPDATE to generate data.

652 **Training** To train an EOGT, we cannot rely on the idea that we can just train a GF and then replace
653 each of its trees by generative trees. To take a concrete example of how this can be a bad idea in the
654 context of missing data imputation, we have observed empirically that a generative forest (GF) can
655 have many trees whose node's observation variables are the *same* within the tree. Taken independently
656 of the forest, such trees would only model marginals, but in a GF, sampling is dependent on the
657 other trees and Lemma 4.4 guarantees the global accuracy of the forest. *However*, if we then replace
658 the GF by an EOGT with the same trees and use them for missing data imputation, this results in
659 imputation at the mode of the marginal(s) (exclusively), which is clearly suboptimal. To avoid this,
660 we have to use a specific training for an EOGT, and to do this, it is enough to change a key part of
661 training in `splitPred`, the computation of probabilities $p_{\mathbf{R}}[.]$ used in (4). Suppose \varkappa very small in
662 Assumption A. To evaluate a split at some leaf, we observe (suppose we have a single tree)

$$p_{\mathbf{R}}[\mathcal{X}_{\lambda_f}] = p_{\mathbf{R}}[\lambda_f] \cdot p_{\mathbf{R}}[\mathcal{X}_{\lambda_f} | \mathcal{X}_{\lambda}] \approx p_{\mathbf{R}}[\lambda_f] \cdot p_{\mathbf{U}}[\mathcal{X}_{\lambda_f} | \mathcal{X}_{\lambda}]$$

663 (and the same holds for λ_t). Extending this to multiple trees and any $\mathcal{C} \in \mathcal{P}(\mathcal{T})$, we get the computation
664 of any $p_{\mathbf{R}}[\mathcal{C}_f]$ and $p_{\mathbf{R}}[\mathcal{C}_t]$ needed to measure the new (4) for a potential split. Crucially, it does not
665 necessitate to split the empirical measure at \mathcal{C} but just relies on computing $p_{\mathbf{U}}[\mathcal{C}_v | \mathcal{C}], v \in \{f, t\}$: with
666 the product measure and since we make axis-parallel splits, it can be done in $O(1)$, thus substantially
667 reducing training time.

668 **More with ensembles of generative trees** we can follow the exact same algorithmic steps as
669 for generative trees to perform missing data imputation and density estimation, but use branching
670 probabilities in the trees to decide branching, instead of relying on the empirical measure at tuples
671 of nodes, which we do not have access to anymore. For this, we rely on the approximation (31)
672 in Lemma C. A key difference with a GF is that no tuple can have zero density because branching
673 probabilities are in $(0, 1)$ in each generative tree.

674 V Appendix on experiments

675 V.1 Domains

676 `ringGauss` is the seminal 2D ring Gaussians appearing in numerous GAN papers [44]; those are
677 eight (8) spherical Gaussians with equal covariance, sampling size and centers located on regularly

Domain	Tag	Source	Missing data ?	m	d	# Cat.	# Num.
iris	-	UCI	No	150	5	1	4
*ringGauss	ring	-	No	1 600	2	-	2
*circGauss	circ	-	No	2 200	2	-	2
*gridGauss	grid	-	No	2 500	2	-	2
forestfires	for	OpenML	No	517	13	2	11
*randGauss	rand	-	No	3 800	2	-	2
tictactoe	tic	UCI	No	958	9	9	-
ionosphere	iono	UCI	No	351	34	2	32
student_performance_mat	stm	UCI	No	396	33	17	16
winered	wred	UCI	No	1 599	12	-	12
student_performance_por	stp	UCI	No	650	33	17	16
analcatdata_supreme	ana	OpenML	No	4 053	8	1	7
abalone	aba	UCI	No	4 177	9	1	8
kc1	-	OpenML	No	2 110	22	1	21
winewhite	wwhi	UCI	No	4 898	12	-	12
sigma-cabs	-	Kaggle	Yes	5 000	13	5	8
compas	comp	OpenML	No	5 278	14	9	5
artificial_characters	arti	OpenML	No	10 218	8	-	8
jungle_chess	jung	OpenML	No	44 819	8	1	6
open-policing-hartford	-	SOP	Yes	18 419	20	16	4
electricity	elec	OpenML	No	45 312	9	2	7

Table A1: Public domains considered in our experiments (m = total number of examples, d = number of features), ordered in increasing $m \times d$. "Cat." is a shorthand for categorical (nominal / ordinal / binary); "Num." stands for numerical (integers / reals). (*) = simulated, URL for OpenML: <https://www.openml.org/search?type=data&sort=runs&id=504&status=active>; SOP = Stanford Open Policing project, <https://openpolicing.stanford.edu/> (see text). "Tag" refers to tag names used in Tables 2 and A6.

678 spaced (2-2 angular distance) and at equal distance from the origin. gridGauss was generated as a
 679 decently hard task from [11]: it consists of 25 2D mixture spherical Gaussians with equal variance
 680 and sampled sizes, put on a regular grid. circGauss is a Gaussian mode surrounded by a circle,
 681 from [44]. randGauss is a substantially harder version of ringGauss with 16 mixture components,
 682 in which covariance, sampling sizes and distances on sightlines from the origin are all random, which
 683 creates very substantial discrepancies between modes.

684 V.2 Algorithms configuration and choice of parameters

685 **GF.BOOST** We have implemented GF.BOOST in Java, following Algorithm 3's blueprint. Our
 686 implementation of tree and leaf in Steps 2.1, 2.2 is simple: we pick the heaviest leaf among all
 687 trees (with respect to R). The search for the best split is exhaustive unless the variable is categorical
 688 with more than a fixed number of distinct modalities (22 in our experiments), above that threshold, we
 689 pick the best split among a random subset. We follow [29]'s experimental setting: in particular, the
 690 input of our algorithm to train a generator is a .csv file containing the training data *without any further*
 691 *information*. Each feature's domain is learned from the training data only; while this could surely and
 692 trivially be replaced by a user-informed domain for improved results (e.g. indicating a proportion's
 693 domain as [0%, 100%], informing the complete list of socio-professional categories, etc.) — and is
 694 in fact standard in some ML packages like weka's ARFF files, we did not pick this option to alleviate
 695 all side information available to the GT learner. Our software automatically recognizes three types of
 696 variables: nominal, integer and floating point represented.

697 **Comments on code** (provided with submission) For the interested reader, we give here some specific
 698 implementation details regarding two classical bottlenecks on tree-based methods (this also applies
 699 to the supervised case of learning decision trees): handling features with large number of modalities
 700 and handling continuous features.
 701

702 We first comment the case where the number of modalities of a feature is large. The details can be
703 found in the code in `File Algorithm.java` (class `Algorithm`), as follows:

- 704 • method `public GenerativeModelBasedOnEnsembleOfTrees learn_geot()` imple-
705 ments both GF.Boost (for generative forests) and its extension to training ensembles of
706 generative trees (Appendix) as a single algorithm.
- 707 • method `Node choose_leaf(String how_to_choose_leaf)` is the method to choose a
708 leaf (Step 2.2). Since it picks the heaviest leaf among all trees, it also implements Step
709 2.1 (we initialize all trees to their roots; having a tree = root for generation does not affect
710 generation).
- 711 • method `public boolean split_generative_forest(Node leaf, HashSet<MeasuredSupportAtTupleOfNodes> tol)` finds the split for a Generative Forest, given
712 a node chosen for the split. Here is how it works:
 - 714 1. It first computes the complete description of possible splits (not their quality yet), using
715 method `public static Vector<FeatureTest> ALL_FEATURE_TESTS(Feature f, Dataset ds)` in `Feature.java`. The method is documented: for continuous
716 features, the list of splits is just a list of evenly spaced splits.
 - 717 2. then it calls `public SplitDetails split_top_down_boosting_generative_forest_fast(Node leaf, boolean [] splittable_feature, FeatureTest [] [] all_feature_tests, HashSet<MeasuredSupportAtTupleOfNodes> tol)`, which returns the best split as follows: (i) it shuffles the potential list
718 of splits and then (ii) picks the best one in the sublist containing the first
719 `Algorithm.MAXIMAL_NUMBER_OF_SPLIT_TESTS_TRIES_PER_BOOSTING_ITERATION`
720 elements (if the list is smaller, the search for the best is exhaustive); this class variable
721 is currently = 1000.

726 Last, we comment how we handle continuous variables: the boosting theory tells us that finding a
727 moderately good split (condition (b) Definition 5.1, main file) is enough to get boosting-compliant
728 convergence (Theorem 5.2, main file), so we have settled for a trivial and efficient mechanism:
729 cut-points are evenly spaced and the number is fixed beforehand. See method `public static`
730 `Vector<FeatureTest> ALL_FEATURE_TESTS(Feature f, Dataset ds)` in `Feature.java`
731 for the details. It turns out that this works very well and shows our theory was indeed concretely used
732 in the design/implementation of the training algorithm.

733 **MICE** We have used the R MICE package V 3.13.0 with two choices of methods for the round robin
734 (column-wise) prediction of missing values: CART [1] and random forests (RF) [40]. In that last case,
735 we have replaced the default number of trees (10) by a larger number (100) to get better results. We
736 use the default number of round-robin iterations (5). We observed that random forests got the best
737 results so, in order not to multiply the experiments reported and perhaps blur the comparisons with
738 our method, we report only MICE's results for random forests.

739 **TENSORFLOW** To learn the additional Random Forests involved in experiments GEN-DISCRIM,
740 we used Tensorflow Decision Forests library[§]. We use 300 trees with max depth 16. Attribute
741 sampling: `sqrt(number attributes)` for classification problems, `number attributes / 3` for regression
742 problems (Breiman rule of thumb); the min #examples per leaf is 5.

743 **CT-GAN** We used the Python implementation[¶] with default values [45].

744 **Adversarial Random Forests** We used the R code of the generator FORGE made available from
745 the paper [42]^{||}, learning forests containing a variable number of trees in {10, 50, 100, 200}. We noted
746 that the code does not run when the dataset has missing values and we also got an error when trying
747 to run the code on kc1.

[§]<https://github.com/google/yggdrasil-decision-forests/blob/main/documentation/learners.md>

[¶]<https://github.com/sdv-dev/CTGAN>

^{||}https://github.com/bips-hb/arf_paper

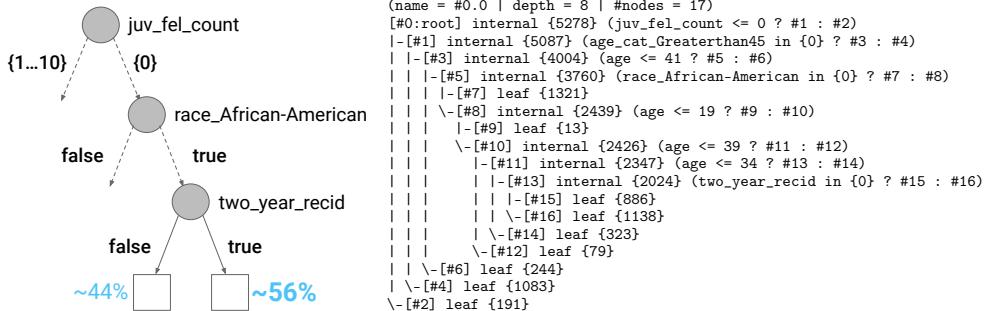


Table A2: SCRUTINIZE: Tree in an example of Generative Forest on sensitive domain `compas`, with 17 nodes and part of its graph sketched following Figure 2's convention, modeling a bias learned from data. Ensemble of such small trees can be *very* accurate: on missing data imputation, they compete with or beat MICE comparatively using 7 000 trees for imputation (see text for details).

748 **Vine Copulas AutoEncoders** We used the Python code available from the paper [39]**, which
 749 processes only fully numerical datasets. We got a few errors when trying to run the code on `compas`.

750 **Kernel Density Estimation** We used the R code available in the package `npudens` with default
 751 values following the approach of [23]. We tried several kernels but ended up with sticking to the
 752 default choices, that seemed to provide overall some of the best results. Compared to us with GF,
 753 KDE's code is extremely compute-intensive as on domains below `tictactoe` in Table A1: it took
 754 orders of magnitude more than ours to get the density values on all folds, so we did not run it on the
 755 biggest domains. Notice that in our case, computation time includes not just training the generative
 756 model, but also, *at each applicable iteration J* in GF.BOOST, the computation of the same number of
 757 density values as for KDE.

758 **Computers used** We ran part of the experiments on a Mac Book Pro 16 Gb RAM w/ 2 GHz
 759 Quad-Core Intel Core i5 processor, and part on a desktop Intel(R) Xeon(R) 3.70GHz with 12 cores
 760 and 64 Gb RAM.

761 V.3 Supplementary results

762 Due to the sheer number of tables to follow, we shall group them according to topics

763 V.V.3.1 Interpreting our models: 'SCRUTINIZE'

764 Table A2 provides experimental results on sensitive real world domain `compas`. It shows that it is
 765 easy to flag potential bias / fairness issues in data by directly estimating conditional probabilities:
 766 considering the GF of this example and ignoring variable age for simplicity, conditionally to having
 767 0 felony count and race being African-American, the probability of generating an observation with
 768 `two_year_recid = true` is .562 (=1138/2024).

769 V.V.3.2 More examples of Table 1 (MF)

770 Tables A3 completes Tables 3 and A2 (main file), displaying that even a stochastic dependence can
 771 be accurately modeled.

772 We provide in Table A4 the density learned on all our four 2D (for easy plotting) simulated domains
 773 and not just `circgauss` as in Table 1 (MF). We see that the observations made in Table 1 (MF)
 774 can be generalized to all domains. Even when the results are less different between 50 stumps in a
 775 generative forest and 1 generative tree with 50 splits for `ringgauss` and `randgauss`, the difference
 776 on `gridgauss` is stark, 1 generative tree merging many of the Gaussians.

**<https://github.com/sdv-dev/Copulas>.

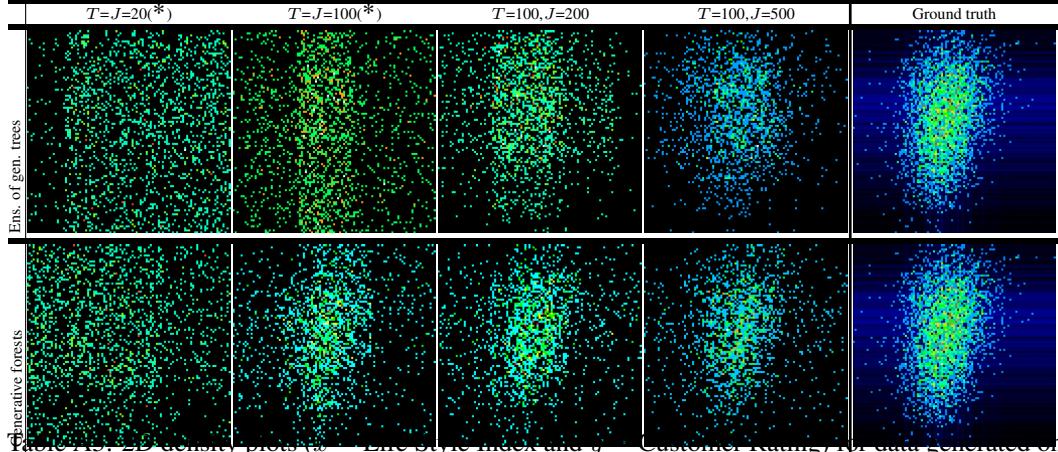


FIGURE A3. 2D density plots of Life Style Index and η Customer Rating, for data generated on sigma_cabs, from models learned with 5% missing features (MCAR), using ensembles of generative trees (top) and generative forests (bottom), for varying total number of trees T and total number of splits J (columns). "*" = all trees are stumps. The rightmost column recalls the domain ground truth for comparison. Each generated dataset contains $m = 2000$ observations.

777 V.V.3.3 The generative forest of Table 1 (MF) developed further

778 One might wonder how a set of stumps gets to accurately fit the domain as the number of stumps
 779 increases. Table A5 provides an answer. From this experiment, we can see that 16 stumps are enough
 780 to get the center mode. The ring shape takes obviously more iterations to represent but still, it takes a
 781 mere few dozen stumps to clearly get an accurate shape, the last iterations just finessing the fitting.

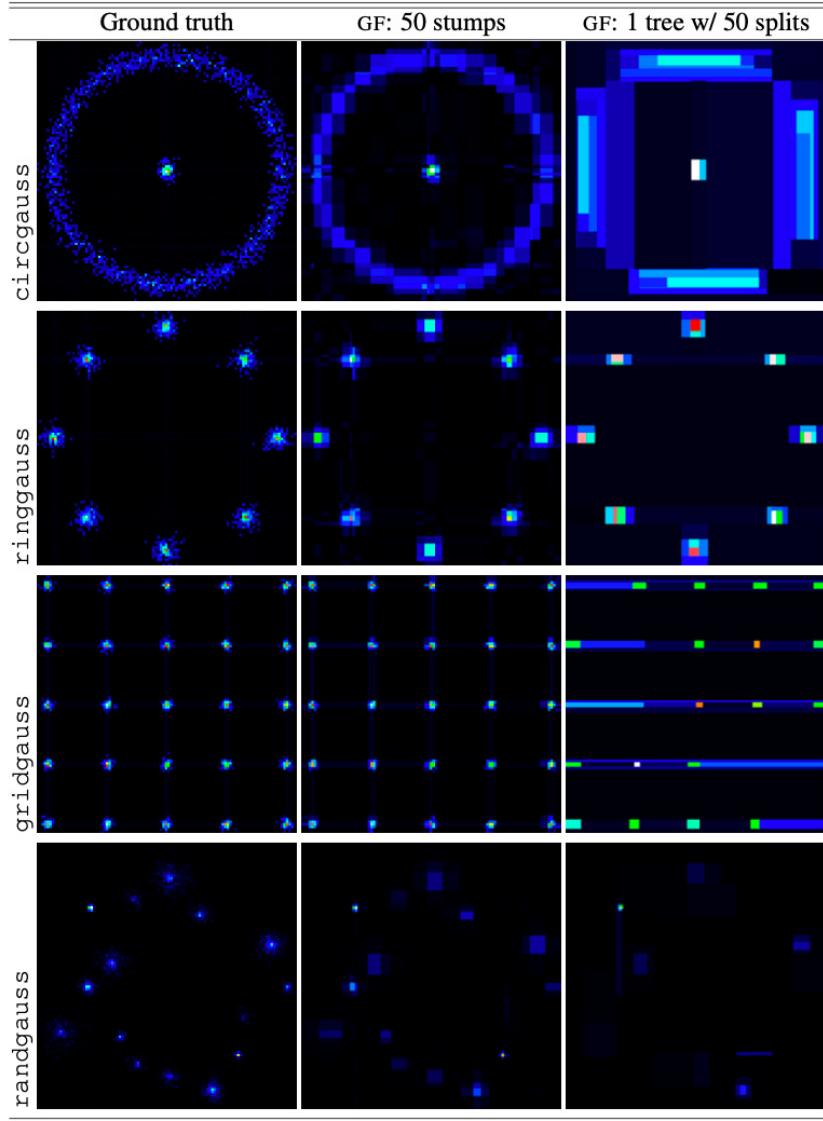


Table A4: Comparison of the density learned by a Generative Forest (GF) consisting of a single tree learned by GF.BOOST with n total splits (*right*) and a set of n GF stumps (1 split per tree) learned by GF.BOOST (*center*). The left picture is the domain, represented using the same heatmap. On each domain, 5% of the data is missing (MCAR = missing completely at random).

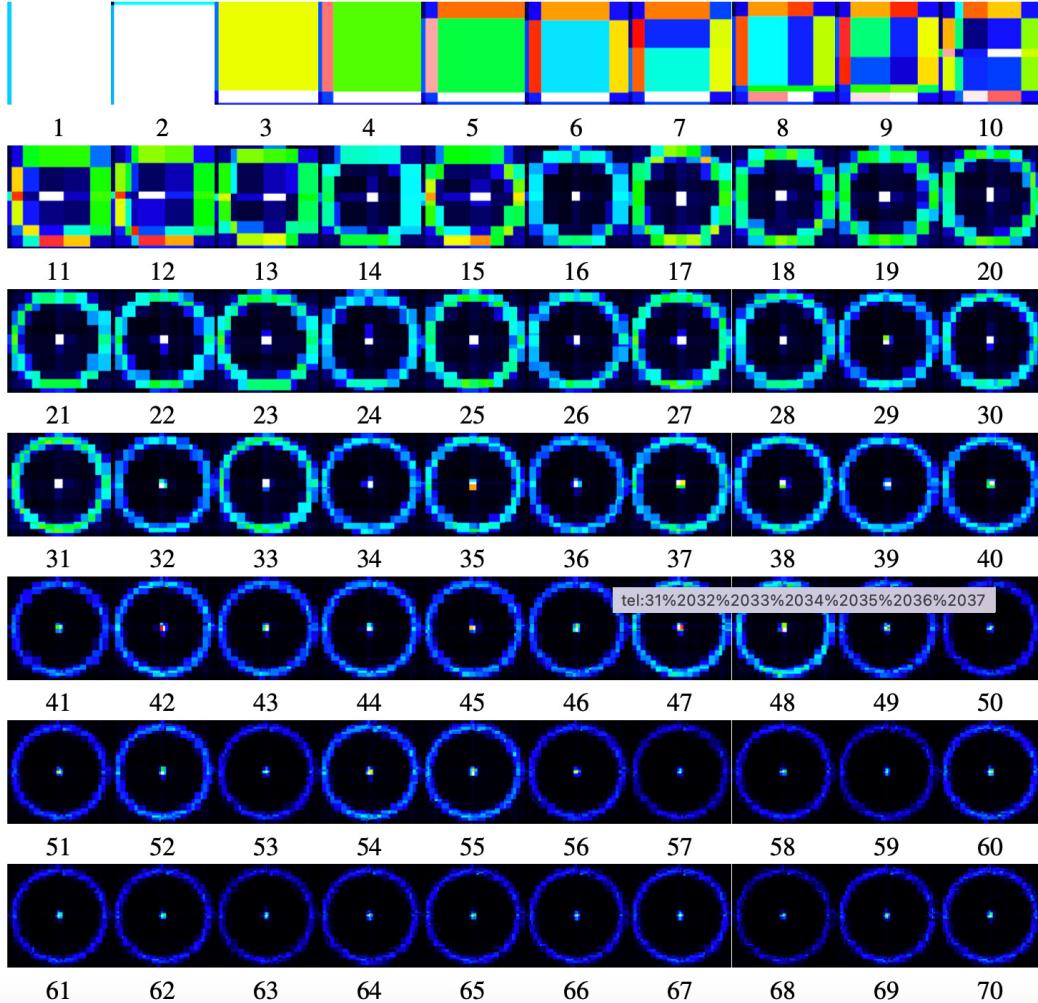


Table A5: Densities learned on `circgauss` by GF.BOOST with generative forests consisting of stumps, for values of $T = J$ in $\{1, 2, \dots, 70\}$. The models quickly capture the ring and the middle dense Gaussian. The boxy-shape of the densities, due to the axis-parallel splits, get less noticeable as $T = J$ exceeds a few dozen.

782 V.V.3.4 Experiment LIFELIKE *in extenso*

783 The objective of the experiment is to evaluate whether a generative model is able to create "realistic"
 784 data. Tabular data is hard to evaluate visually, unlike *e.g.* images or text, so we have considered
 785 a simple evaluation pipeline: we create for each domain a 5-fold stratified experiment. After a
 786 generative model has been trained, we generate the same number of observations as in the test fold
 787 and compute the optimal transport (OT) distance between the generated sample and the fold's test
 788 sample. To fasten the computation of OT costs, we use Sinkhorn's algorithm [8] with an ε -entropic
 789 regulariser, for some $\varepsilon = 0.5$ which we observed was the smallest in our experiments to run with all
 790 domains without leading to numerical instabilities. To balance the importance of categorical features
 791 (for which the cost is binary, depending on whether the guess is right or not) and numerical features,
 792 we normalize numerical features with the domain's mean and standard deviation prior to computing
 793 OT costs. We then compare our method to three possible contenders: Adversarial Random Forests
 794 (ARF) [42], CT-GANs [45] and Vine copula autoencoders (VCAE) [39]. All these approaches rely
 795 on models that are very different from each other. Adversarial Random Forests (ARF) are tree-based
 796 generators that work as follows to generate one observation: one first sample uniformly at random a
 797 tree, then samples a leaf in the tree. The leaf is attached to a distribution which is then used to sample
 798 the observation. As we already pointed out in Section 2, there are two main differences with our

799 models. From the standpoint of the model's distribution, if one takes the (non-empty) support from a
800 tuple of leaves, its probability is obtained from a weighted average of the tree's distributions in [42]
801 while it is obtained from a product of theirs in our case. Hence, at similar tree sizes, the modelling
802 capability of the set of trees tips in our favor, but it is counterbalanced by the fact that their approach
803 equip leaves with potentially complex distributions (e.g. truncated Gaussians) whereas we stick to
804 uniform distributions at the leaves. A consequence of the modeling of ARFs is that each tree has to
805 separately code for a good generator: hence, it has to be big enough or the distributions used at the
806 leaves have to be complex enough. In our case, as we already pointed out, a small number of trees
807 (sometimes even *stumps*) can be enough to get a fairly good generator, even on real-world domains.
808 Our second contender, CT-GAN [45], fully relies on neural networks so it is *de facto* substantially
809 different from ours. The last, VCAE [39] relies on a subtle combination of deep nets and graphical
810 models learned in the latent space. ARFs are trained with a number of trees $T \in \{10, 50, 100, 200\}$.
811 ARFs include an algorithm to select the tree size so we do not have to select it. CT-GANs are trained
812 with a number of epochs $E \in \{10, 100, 300, 1000\}$. We tested VCAE with all types of vine copulas
813 (center, direct and regular).

814 We compare those models with generative forests with $T = 500$ trees and trained for a total of
815 $J = 2000$ iterations, for all domains considered. Compared to the trees learned by ARFs, the total
816 number of splits we use can still be small compared to theirs, in particular when they learn $T \geq 100$
817 trees. Each table against ARFs, CT-GANs and VCAE has a different size and contains only a subset
818 of the whole 21 domains: VCAE can only process fully numerical domains, ARFs do not process
819 domains with missing values and we also got issues running contenders on several domains: those
820 are discussed in Appendix, Section V.2.

821 **Results vs Adversarial Random Forests** Table 2 (main file) presents the results obtained against
822 adversarial random forests. A first observation, not visible in the table, is that ARFs indeed tend to
823 learn big models, typically with dozens of nodes in each tree. For the largest ARFs with 100 or 200
824 trees, this means in general a total of thousands of nodes in models. A consequence, also observed
825 experimentally, is that there is sometimes little difference in performance in general between models
826 with a different number of trees in ARFs as each tree is in fact an already good generator. In our case,
827 this is obviously not the case. Noticeably, the performance of ARFs is not monotonic in the number
828 of trees, so there could be a way to find the appropriate number of trees in ARFs to buy the slight
829 (but sometimes significant) increase in performance in a domain-dependent way. In our case, there is
830 obviously a dependency in the number of trees chosen as well. From the standpoint of the optimal
831 transport metric, even when ARFs make use of distributions at their tree leaves that are much more
832 powerful than ours and in fact fit to some of our simulated domains (Gaussians used in `circgauss`,
833 `randgauss` and `gridgauss`), we still manage to compete or beat ARFs on those simulated domains.

834 Globally, we manage to beat ARFs on almost all runs, and very significantly on many of them.
835 Since training generative forests does not include a mechanism to select the size of models, we
836 have completed this experiment by another one on which we learn much smaller generative forests.
837 The corresponding Table is in Appendix, Section V.V.3.4. Because we clearly beat ARFs on
838 `student_performance_mat` and `student_performance_por` in Table 2 but are beaten by ARFs
839 for much smaller generative forests, we conclude that there could exist mechanisms to compute the
840 "right size" of our models, or to prune big models to get models with the right size. Globally however,
841 even with such smaller models, we still manage to beat ARFs on a large majority of cases, which is a
842 good figure given the difference in model sizes. The ratio "results quality over model size" tips in our
843 favor and demonstrates the potential in using all trees in a forest to generate an observation (us) vs
844 using a single of its trees (ARFs), see Figure 1.

845 **Results vs CT-GANs** Table A6 summarizes our results against CT-GAN using various numbers
846 of training epochs (E). In our case, our setting is the same as versus ARFs: we stick to $T = 500$
847 trees trained for a total number of $J = 2000$ iterations in the generative forest, for all domains. We
848 see that generative forests consistently and significantly outperform CT-GAN on nearly all cases.
849 Furthermore, while CT-GANs performances tend to improve with the number of epochs, we observe
850 that on a majority of datasets, performances at 1 000 epochs are still far from those of generative
851 forests and already induce training times that are far bigger than for generative forests: for example,
852 more than 6 hours per fold for `stm` while it takes just a few seconds to train a generative forest. Just
853 like we did for ARFs, we repeated the experiment vs CT-GAN using much smaller generative forests
854 ($T = 200, J = 500$). The results are presented in Appendix, Section V.V.3.4. There is no change in

domain tag	us (GF) $T=500, J=2\,000$	CT-GANs							
		$E = 10$	$p\text{-val}$	$E = 100$	$p\text{-val}$	$E = 300$	$p\text{-val}$	$E = 1000$	$p\text{-val}$
ring	0.285 \pm 0.008	0.457 \pm 0.058	0.0032	0.546 \pm 0.079	0.0018	0.405 \pm 0.044	0.0049	0.351 \pm 0.042	0.0183
circ	0.351 \pm 0.005	0.848 \pm 0.300	0.0213	0.480 \pm 0.040	0.0019	0.443 \pm 0.014	ε	0.435 \pm 0.075	0.0711
grid	0.390 \pm 0.002	0.749 \pm 0.417	0.1268	0.459 \pm 0.031	0.0085	0.426 \pm 0.022	0.0271	0.408 \pm 0.019	0.1000
for	1.108 \pm 0.105	9.796 \pm 5.454	0.0049	1.899 \pm 0.289	0.0045	1.532 \pm 0.178	0.0244	1.520 \pm 0.311	0.0410
rand	0.286 \pm 0.003	0.746 \pm 0.236	0.0119	0.491 \pm 0.063	0.0021	0.368 \pm 0.015	0.0002	0.327 \pm 0.024	0.0288
tic	0.575 \pm 0.002	0.601 \pm 0.003	0.0002	0.581 \pm 0.001	0.0082	0.586 \pm 0.006	0.0294	0.584 \pm 0.003	0.0198
iono	1.167 \pm 0.074	2.263 \pm 0.035	ε	2.084 \pm 0.052	ε	1.984 \pm 0.136	ε	1.758 \pm 0.137	0.0002
stm	0.980 \pm 0.032	1.511 \pm 0.074	ε	1.189 \pm 0.045	0.0002	1.168 \pm 0.054	0.0056	1.167 \pm 0.041	0.0013
wred	0.980 \pm 0.032	2.836 \pm 0.703	0.0044	2.004 \pm 0.090	ε	1.825 \pm 0.127	0.0001	1.384 \pm 0.047	ε
stp	0.976 \pm 0.018	1.660 \pm 0.161	0.0006	1.141 \pm 0.018	0.0001	1.206 \pm 0.046	0.0007	1.186 \pm 0.052	0.0019
ana	0.368 \pm 0.014	0.542 \pm 0.135	0.0050	0.439 \pm 0.056	0.0202	0.404 \pm 0.012	0.0397	0.436 \pm 0.036	0.0051
aba	0.490 \pm 0.028	1.689 \pm 0.053	ε	1.463 \pm 0.094	ε	0.754 \pm 0.040	0.0009	0.657 \pm 0.026	0.0006
wwhi	1.064 \pm 0.003	1.770 \pm 0.160	0.0005	1.849 \pm 0.064	ε	1.284 \pm 0.029	ε	1.158 \pm 0.009	ε
arti*	0.821 \pm 0.004	1.388 \pm 0.109 ₄	0.0020	1.209 \pm 0.018 ₄	0.0005	1.026 \pm 0.013 ₄	0.0001	0.959 \pm 0.013 ₃	0.0038
wins / lose for us \rightarrow	14 / 0		14 / 0		14 / 0		14 / 0		14 / 0

Table A6: Comparison of results with respect to CT-GANs [45] trained with different number of epochs E . On some domains, indicated with a “*”, CT-GANs crashed on some folds. For such domains, we indicate in index to CT-GANs results the number of folds (out of 5) for which this did *not* happen. In such cases, we restricted the statistical comparisons with us to the folds for which CT-GANs did not crash: the statistical tests take this into account; instead of providing all corresponding average performances for us, we keep giving the average performance for us on all five folds (leftmost column), which is thus indicative. Other conventions follow Table 2.

855 the conclusion: even with such small models, we still beat CT-GANs, regardless of the number of
 856 training epochs, on all cases (and very significantly on almost all of them).

857 **Results vs VCAE** Table A7 summarizes our results against Gaussian multivariate modelling and
 858 Vine copulas autoencoders [39]^{††}, with all three vine copulas options tested (center, direct and
 859 regular). The setting for our method is the same as vs ARFs and CT-GANs ($T = 500$ trees, $J = 2000$
 860 total iterations). Because those techniques we evaluate our method against can only process fully
 861 numerical datasets, the Table collects all numerical domains from Tables 2 and A6 (when some
 862 variables are binary like in `iono`, we can treat them as numerical). We also indicate the result
 863 of a standard multivariate Gaussian, which can serve as a baseline for us and VCAE. The results
 864 display that pretty much like vs ARFs and CT-GANs, we also consistently outperform Vine copula
 865 approaches, and by a substantial margin in most cases. What is very interesting is that our simulated
 866 domains (`ring`, `circ`, `grid` and `rand`) are all Gaussian mixtures, and it appears that we also beat
 867 consistently Gaussian multivariate approaches. Our models do not have a representational bias that
 868 would make them especially fit for Gaussians, so we take this result as a good experimental validation
 869 of the general goodness-of-fit of our models and training algorithm for the data generation task. Just
 870 like we did for ARFs, we repeated the experiment vs VCAE using much smaller generative forests
 871 ($T = 200, J = 500$). The results are presented in Appendix, Section V.V.3.4. They support the same
 872 conclusion as for the bigger models as we still consistently beat VCAE on all domains.

873 **More results with small generative forests** Tables A8, A9 and A10 present the results of generative
 874 forests vs adversarial random forests, CT-GANs, Gaussian multivariate modelling and Vine copulas
 875 (those two last approaches carried out only for fully numerical domains in Table A8), when the size
 876 of our models is substantially smaller than in the main file ($T = 200, J = 500$). We still manage
 877 to compete or beat ARFs on simulated domains for which modelling the leaf distributions in ARFs
 878 should represent an advantage (Gaussians used in `circgauss`, `randgauss` and `gridgauss`), even
 879 more using comparatively very small models compared to ARFs. We believe this could be due to
 880 two factors: (i) the fact that in our models all trees are used to generate each observation (which
 881 surely facilitates learning small and accurate models) and (ii) the fact that we do not generate data
 882 during training but base our splitting criterion on an *exact* computation of the reduction of Bayes
 883 risk in growing trees. If we compute aggregated statistics comparing, for each number of trees in
 884 ARFs, the number of times we win / lose versus ARFs, either considering only significant p -values or
 885 disregarding them, our generative forests always beat ARFs on a majority of domains. Globally, this
 886 means we win on 34 out of 52 cases. The results vs CT-GANs are of the same look and feel as in the
 887 main file with larger generative forests: generative forests beat them in all cases, and very significantly

^{††}<https://github.com/sdv-dev/Copulas>.

domain tag	us (GF) $J=500, T=2\,000$	Multivariate Gaussian		$p\text{-val}$	Vine copulas		direct		$p\text{-val}$	regular		$p\text{-val}$
		center	$p\text{-val}$		direct	$p\text{-val}$	regular	$p\text{-val}$		regular	$p\text{-val}$	
ring	0.285 \pm 0.008	0.392 \pm 0.006	ϵ	0.394 \pm 0.016	ϵ	0.391 \pm 0.012	ϵ	0.388 \pm 0.005	ϵ			
circ	0.351 \pm 0.005	0.415 \pm 0.013	0.0001	0.409 \pm 0.005	ϵ	0.411 \pm 0.004	ϵ	0.403 \pm 0.003	0.0001			
grid	0.390 \pm 0.002	0.400 \pm 0.003	0.0045	0.400 \pm 0.005	0.0187	0.399 \pm 0.002	ϵ	0.400 \pm 0.003	0.0066			
rand	0.286 \pm 0.003	0.414 \pm 0.003	ϵ	0.415 \pm 0.005	ϵ	0.414 \pm 0.003	ϵ	0.418 \pm 0.010	ϵ			
iono	1.167 \pm 0.074	1.928 \pm 0.553	0.0279	1.836 \pm 0.153	0.0002	1.807 \pm 0.082	0.0004	1.722 \pm 0.041	0.0001			
wred	0.980 \pm 0.032	1.106 \pm 0.043	0.0002	1.544 \pm 0.024	ϵ	1.383 \pm 0.038	ϵ	1.365 \pm 0.099	0.0004			
wwhi	1.064 \pm 0.003	1.181 \pm 0.020	0.0001	1.437 \pm 0.049	ϵ	1.484 \pm 0.056	ϵ	1.354 \pm 0.039	ϵ			
comp*	0.532 \pm 0.008	0.960 \pm 0.012 ₅	ϵ	1.004 \pm 0.022 ₅	ϵ	0.988 \pm 0.022 ₄	ϵ	No result				
arti	0.830 \pm 0.007	0.984 \pm 0.019	ϵ	1.068 \pm 0.016	ϵ	1.154 \pm 0.010	ϵ	1.102 \pm 0.030	ϵ			
wins / lose for us \rightarrow	9 / 0			9 / 0		9 / 0		8 / 0				

Table A7: Comparison of results with respect to Gaussian multivariate modelling and Vine copulas autoencoders (VCAE) with the three types of copulas indicated. See Table A6 for the meaning of the "*" appearing in some domains. "No results" mean Vine did not provide results. Other conventions follow Table 2 (only numerical datasets processed).

domain tag	us (GF) $T=200, J=500$	Adversarial Random Forests (ARFs)							
		$T = 10$	$p\text{-val}$	$T = 50$	$p\text{-val}$	$T = 100$	$p\text{-val}$	$T = 200$	$p\text{-val}$
iris	0.529 \pm 0.099	0.530 \pm 0.064	0.9890	0.517 \pm 0.081	0.8645	0.466 \pm 0.050	0.1099	0.501 \pm 0.043	0.6548
ring	0.283 \pm 0.005	0.289 \pm 0.006	0.0651	0.286 \pm 0.007	0.4903	0.288 \pm 0.010	0.4292	0.286 \pm 0.007	0.7442
circ	0.356 \pm 0.002	0.354 \pm 0.002	0.4561	0.356 \pm 0.008	0.7882	0.350 \pm 0.002	0.0591	0.355 \pm 0.005	0.7446
grid	0.392 \pm 0.002	0.394 \pm 0.003	0.4108	0.391 \pm 0.002	0.3403	0.392 \pm 0.001	0.8370	0.394 \pm 0.002	0.3423
for	1.086 \pm 0.124	1.272 \pm 0.281	0.1420	1.431 \pm 0.337	0.1142	1.311 \pm 0.255	0.1215	1.268 \pm 0.209	0.0370
rand	0.290 \pm 0.007	0.286 \pm 0.003	0.4488	0.290 \pm 0.005	0.9553	0.289 \pm 0.005	0.9004	0.288 \pm 0.002	0.6542
tic	0.574 \pm 0.001	0.577 \pm 0.001	0.0394	0.578 \pm 0.001	0.0325	0.577 \pm 0.001	0.0657	0.577 \pm 0.001	0.0914
iono	1.203 \pm 0.068	1.431 \pm 0.043	0.7307	1.469 \pm 0.079	0.3402	1.431 \pm 0.058	0.7046	1.420 \pm 0.056	0.8737
stm	1.044 \pm 0.013	1.010 \pm 0.015	0.0080	1.015 \pm 0.018	0.0094	1.001 \pm 0.010	0.0037	1.014 \pm 0.021	0.0329
wred	1.027 \pm 0.037	1.086 \pm 0.036	0.0327	1.072 \pm 0.055	0.1432	1.086 \pm 0.030	0.0095	1.099 \pm 0.031	0.0044
stp	1.029 \pm 0.027	0.999 \pm 0.011	0.0826	1.012 \pm 0.010	0.3421	1.006 \pm 0.011	0.0542	1.003 \pm 0.003	0.0764
ana	0.367 \pm 0.018	0.370 \pm 0.012	0.8229	0.382 \pm 0.015	0.2318	0.368 \pm 0.006	0.9264	0.369 \pm 0.007	0.8722
aba	0.476 \pm 0.017	0.505 \pm 0.046	0.2264	0.484 \pm 0.024	0.2162	0.503 \pm 0.035	0.1847	0.481 \pm 0.031	0.7468
wwhi	1.120 \pm 0.013	1.159 \pm 0.011	0.0013	1.159 \pm 0.006	0.0014	1.170 \pm 0.014	0.0026	1.150 \pm 0.021	0.0327
comp	0.538 \pm 0.015	0.531 \pm 0.012	0.5639	0.534 \pm 0.009	0.5706	0.551 \pm 0.023	0.4057	0.535 \pm 0.033	0.9017
arti	0.830 \pm 0.007	0.849 \pm 0.004	0.0078	0.847 \pm 0.010	0.0012	0.843 \pm 0.009	0.0003	0.849 \pm 0.010	0.0010
jung	0.928 \pm 0.001	0.929 \pm 0.002	0.1720	0.930 \pm 0.002	0.0646	0.929 \pm 0.002	0.2675	0.930 \pm 0.001	0.0081
elec	1.503 \pm 0.035	1.577 \pm 0.088	0.1222	1.543 \pm 0.017	0.0880	1.552 \pm 0.039	0.1228	1.581 \pm 0.025	0.0027
wins / lose for us \rightarrow	13 / 5			11 / 5		12 / 5		12 / 6	

Table A8: Comparison of results for ARFs [42] with different number J of trees, and Generative Forests (us, GF) where the number of trees and number of iterations are **fixed** ($T = 200$ trees, total of $J = 500$ splits in trees). Values are shown on the form average over the 5-folds \pm standard deviation. The p -values for the comparison of our results ("us") and ARFs are shown. Values appearing in green in p -vals mean (i) the average regularised OT cost for us (GF) is smaller than that of ARFs and (ii) $p < 0.1$, meaning we can conclude that our method outperforms ARFs. Values appearing in magenta in p -vals mean (i) the average regularised OT cost for us (GF) is larger than that of ARFs and (ii) $p < 0.1$, meaning we can conclude that our method is outperformed by ARFs. The last row summarizes the number of times we win / lose against ARFs according to the average metric used. To reduce size, domains are named by a tag: the correspondence with domains and their characteristics is in Appendix, Table A1.

888 in most of them. Finally, the results vs Vine copulas and Gaussian multivariate modelling, show the
889 same trend as for bigger models as we beat all approaches, albeit with a slightly smaller number of
890 statistically significant wins.

891 V.V.3.5 Comparison with "the optimal generator": GEN-DISCRIM

892 Rather than compare our method with another one, the question we ask here is "can we generate data
893 that looks like domain's data"? We replicate the experimental pipeline of [29]. In short, we shuffle
894 a 3-partition (say $\mathbf{R}_1, \mathbf{R}_2, \mathbf{R}_3$) of the training data in a $3! = 6$ -fold CV, then train a generator \mathbf{G}_1
895 from \mathbf{R}_1 (we use generative forests with GF.BOOST), then train a random forest (RF) to distinguish
896 between \mathbf{G}_1 and \mathbf{R}_2 , and finally estimate its test accuracy on \mathbf{G}_1 vs \mathbf{R}_3 . The *lower* this accuracy, the
897 less distinguishable are fakes from real and thus the *better* the generator. We consider 3 competing

domain tag	us (GF) $T=200, J=500$	CT-GANs							
		$E = 10$	$p\text{-val}$	$E = 100$	$p\text{-val}$	$E = 300$	$p\text{-val}$	$E = 1000$	$p\text{-val}$
ring	0.283 \pm 0.005	0.457 \pm 0.058	0.0029	0.546 \pm 0.079	0.0020	0.405 \pm 0.044	0.0048	0.351 \pm 0.042	0.0299
circ	0.356 \pm 0.002	0.848 \pm 0.300	0.0217	0.480 \pm 0.040	0.0024	0.443 \pm 0.014	0.0001	0.435 \pm 0.075	0.0816
grid	0.392 \pm 0.002	0.749 \pm 0.417	0.1279	0.459 \pm 0.031	0.0096	0.426 \pm 0.022	0.0233	0.408 \pm 0.019	0.1239
for	1.086 \pm 0.124	9.796 \pm 5.454	0.0241	1.899 \pm 0.289	0.0028	1.532 \pm 0.178	0.0020	1.520 \pm 0.311	0.0313
rand	0.290 \pm 0.007	0.746 \pm 0.236	0.0129	0.491 \pm 0.063	0.0017	0.368 \pm 0.015	0.0003	0.327 \pm 0.024	0.0124
tic	0.574 \pm 0.001	0.601 \pm 0.003	ε	0.581 \pm 0.001	0.0024	0.586 \pm 0.006	0.0299	0.584 \pm 0.003	0.0228
iono	1.203 \pm 0.068	2.263 \pm 0.035	ε	2.084 \pm 0.052	ε	1.98 \pm 0.136	0.0001	1.758 \pm 0.137	0.0005
stm	1.044 \pm 0.013	1.511 \pm 0.074	0.0001	1.189 \pm 0.045	0.0106	1.168 \pm 0.054	0.0346	1.167 \pm 0.041	0.0036
wred	1.027 \pm 0.037	2.836 \pm 0.703	0.0049	2.004 \pm 0.090	ε	1.825 \pm 0.127	ε	1.384 \pm 0.047	0.0002
stp	1.029 \pm 0.027	1.660 \pm 0.161	0.0010	1.141 \pm 0.018	0.0004	1.206 \pm 0.046	0.0015	1.186 \pm 0.052	0.0057
ana	0.367 \pm 0.018	0.542 \pm 0.135	0.0399	0.439 \pm 0.056	0.0196	0.404 \pm 0.012	0.0231	0.436 \pm 0.036	0.0018
aba	0.476 \pm 0.017	1.689 \pm 0.053	ε	1.463 \pm 0.094	ε	0.754 \pm 0.040	ε	0.657 \pm 0.026	ε
wwhi	1.120 \pm 0.013	1.770 \pm 0.160	0.0009	1.849 \pm 0.064	ε	1.284 \pm 0.029	0.0006	1.158 \pm 0.009	0.0117
arti*	0.830 \pm 0.007	1.388 \pm 0.109 ₄	0.0021	1.209 \pm 0.018 ₄	ε	1.026 \pm 0.013 ₄	0.0002	0.959 \pm 0.013 ₃	0.0051
wins / lose for us \rightarrow	14 / 0			14 / 0		14 / 0		14 / 0	

Table A9: Comparison of results with respect to CT-GANs [45] trained with different number of epochs E . ε means $p\text{-val} < 10^{-4}$. As already mentioned in Table A6 (main file), in some domains, indicated with a "*", CT-GANs crashed on some folds. For such domains, we indicate in index to CT-GANs results the number of folds (out of 5) for which this did *not* happen. In such cases, we restricted the statistical comparisons with us to the folds for which CT-GANs did not crash: the statistical tests take this into account; instead of providing all corresponding average performances for us, we keep giving the average performance for us on all five folds (leftmost column), which is thus indicative. Other conventions follow Table A8.

domain tag	us (GF) $T=200, J=500$	Multivariate Gaussian		$p\text{-val}$	Vine copulas				$p\text{-val}$
		center	$p\text{-val}$		direct	$p\text{-val}$	regular	$p\text{-val}$	
ring	0.283 \pm 0.005	0.392 \pm 0.006	ε	0.394 \pm 0.016	0.0001	0.391 \pm 0.012	ε	0.388 \pm 0.005	ε
circ	0.356 \pm 0.002	0.415 \pm 0.013	0.0002	0.409 \pm 0.005	ε	0.411 \pm 0.004	ε	0.403 \pm 0.003	ε
grid	0.392 \pm 0.002	0.401 \pm 0.003	0.0003	0.400 \pm 0.005	0.0285	0.399 \pm 0.002	0.0095	0.400 \pm 0.003	0.0168
rand	0.290 \pm 0.007	0.414 \pm 0.003	ε	0.415 \pm 0.005	ε	0.414 \pm 0.003	ε	0.418 \pm 0.010	ε
iono	1.203 \pm 0.068	1.928 \pm 0.553	0.0371	1.836 \pm 0.153	0.0001	1.807 \pm 0.082	0.0001	1.722 \pm 0.041	0.0002
wred	1.027 \pm 0.037	1.106 \pm 0.043	0.0154	1.544 \pm 0.024	ε	1.383 \pm 0.038	0.0002	1.365 \pm 0.099	0.0005
wwhi	1.120 \pm 0.013	1.181 \pm 0.020	0.0082	1.437 \pm 0.049	0.0001	1.484 \pm 0.056	0.0003	1.354 \pm 0.039	0.0002
comp*	0.538 \pm 0.015	0.958 \pm 0.012 ₅	ε	1.004 \pm 0.022 ₅	ε	0.988 \pm 0.022 ₄	ε	No result	
arti	0.830 \pm 0.007	0.984 \pm 0.019	ε	1.154 \pm 0.010	ε	1.068 \pm 0.016	ε	1.102 \pm 0.030	ε
wins / lose for us \rightarrow	9 / 0			9 / 0		9 / 0		8 / 0	

Table A10: Comparison of results with respect to Gaussian multivariate modelling and VCAE with the three types of copulas indicated. ε means $p\text{-val} < 10^{-4}$. See Table A9 for the meaning of the "*" appearing in some domains. "No results" mean Vine did not provide results. Other conventions follow Table A9 (only numerical datasets processed).

898 baselines to our models: (i) the "optimal" one, COPY, which consists in replacing G_1 by a set of
899 real data, (ii) the "worst" one, UNIF(orm), which uniformly samples data, and (iii) GF.BOOST for
900 $T = 1$, which learns a generative tree (GT). We note that this pipeline is radically different from
901 the one used in [42]. In this pipeline, domains used are supervised (the data includes a variable to
902 predict). A generator is trained from some training data, then generates an amount of data equal to
903 the size of the training data. Then, two classifiers are trained, one from the original training data, one
904 from the generated data, to predict for the variable to predict and their performances are compared on
905 this basis (the higher the accuracy, the better). There is a risk in this pipeline that we have tried to
906 mitigate with our sophisticated pipeline: if the generator consists just in copying its training data, it is
907 guaranteed to perform well according to [42]'s metric. Obviously, this would not be a good generator
908 however. Our pipeline would typically prevent this, R_2 being different from R_3 .
909 **Results** Tables A11 and A12 provide the results we got, including statistical p -values for the test of
910 comparing our method to COPY. They display that it is possible to get GFs generating realistically
911 looking data: for `iris`, the Student t -tests show that we can keep H_0 = "GFs perform identically
912 to COPY" for $J = T = 40$ ($p > .2$). For `mat`, the result is quite remarkable not because of the
913 highest p -val, which is not negligible (> 0.001), but because to get it, it suffices to build a GF in
914 which the total number of feature occurrences ($J = 80$) is barely three times as big as the domain's
915 dimension ($d = 33$). The full tables display a clear incentive to grow further the GFs as domains
916 increase in size. However, our results show, pretty much like the experiment against Adversarial

917 Random Forests, that there could be substantial value of learning the right model size: `iris` and
918 `student_performance_por` clearly show nontrivial peaks of p -values, for which we would get the
919 best models compared to COPY.

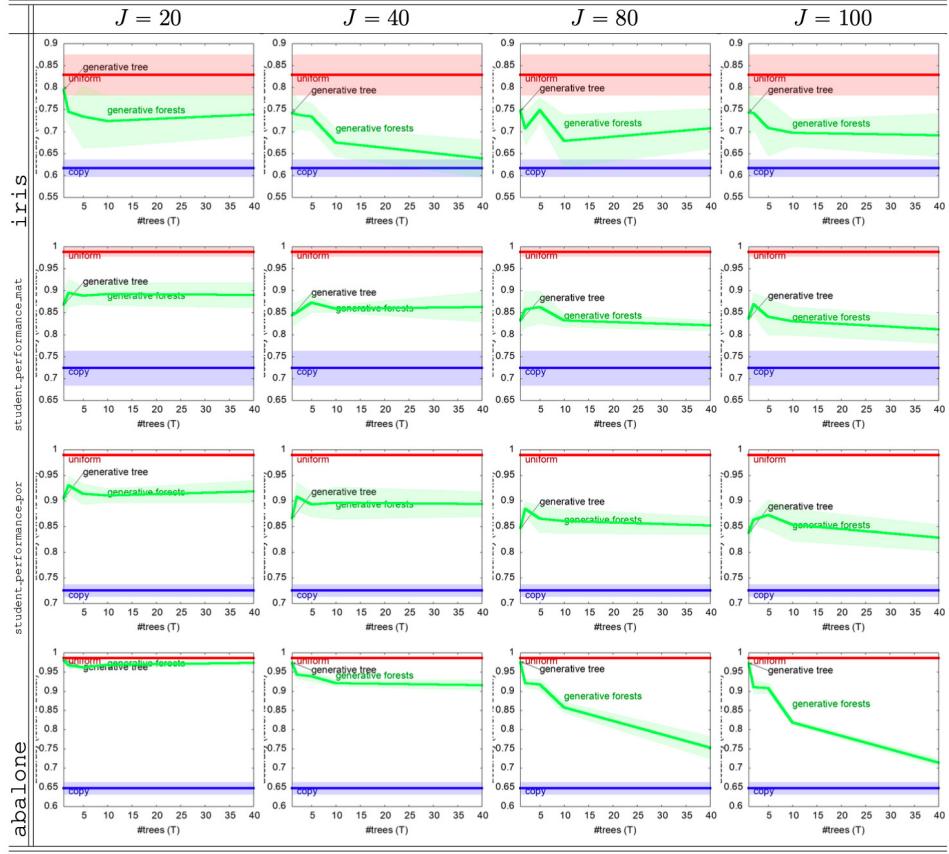


Table A11: Results of generative forests on GEN-DISCRIM for four UCI domains, ordered, **from top to bottom, in increasing domain size**. In each row, the accuracy of distinguishing fakes from real is plotted (the *lower*, the *better* the technique) for four contenders: UNIFORM ("generates observation uniformly at random"), COPY ("optimal" contender using real data as generated), GF.boost inducing generative forests and generative trees (indicated for $T = 1$). A clear pattern emerges, that as the domain size grows, increasing J buys improvements and, if J is large enough, increasing T also improves results. See Table A12 for comparisons (us vs COPY) in terms of p -values.

920 V.V.3.6 Full comparisons with MICE on missing data imputation

921 The main file provides just two examples of domains for the comparison, `abalone` and
 922 `analcatdata_supreme`. We provide here more complete results on the domain used in the main
 923 file and results on more domains in the form of one table for each additional domain:

- 924 • Table A13: experiments on `analcatdata_supreme` completing the results shown in Table
 925 5 (MF);
- 926 • Table A14: experiments on `abalone` completing the results shown in Table 5 (MF);
 927 (tables below are ordered in increasing domain size, see Table A1)
- 928 • Table A15: experiments on `iris`;
- 929 • Table A16: experiments on `ringgauss`;
- 930 • Table A17: experiments on `circgauss`;
- 931 • Table A18: experiments on `gridgauss`;
- 932 • Table A19: experiments on `randgauss`;
- 933 • Table A20: experiments on `student_performance_mat`;
- 934 • Table A21: experiments on `student_performance_por`;

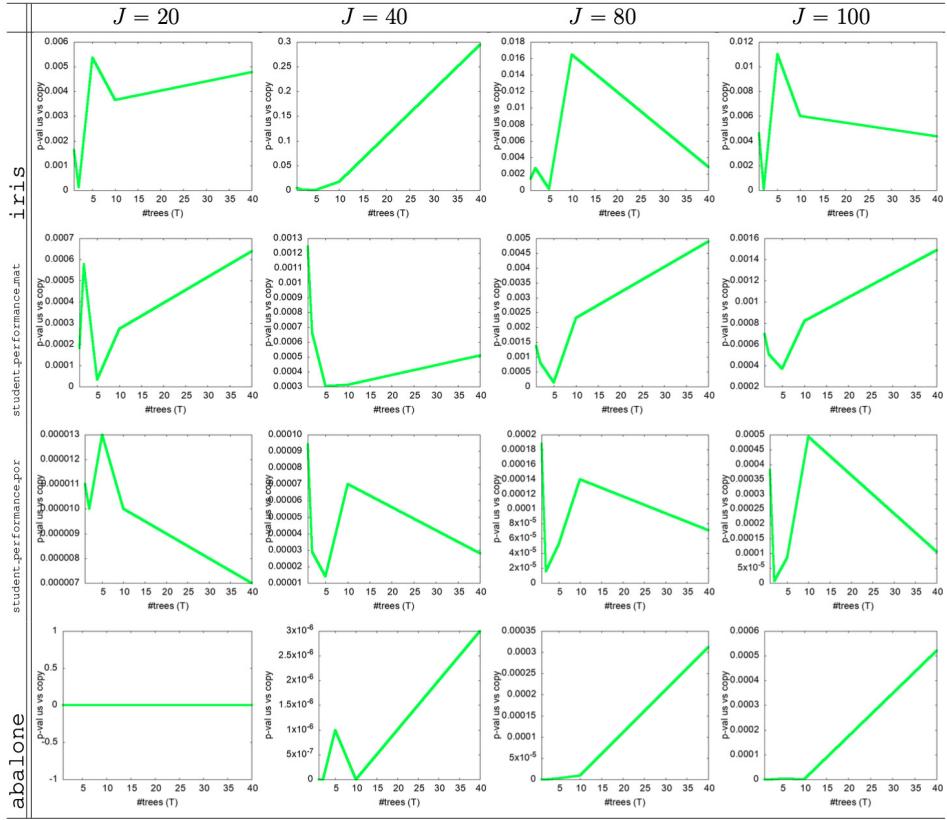


Table A12: p -values for Student tests with H_0 being "GF and COPY perform identically on GEN-DISCRIM" (for the domains in Table A11). Large values indicate we can keep H_0 and thus consider that GF is successful at generating "realistic" data. Clearly, GF are successful for *iris* (p increases with T up to $p \sim 0.3$ for $J = 40$). For the other domains, Table A11 suggests increasing J or T . The curves for *abalone* reveal a dramatic relative increase of p as (J, T) increases, starting from $p \sim 0, \forall T$ for $J = 20$.

- 935 • Table A22: experiments on *kc1*;
 936 • Table A23: experiments on *sigma_cabs*;
 937 • Table A24: experiments on *compas*;
 938 • Table A25: experiments on *open_policing_hartford*;
- 939 The large amount of pictures to process may make it difficult to understand at first glance how our
 940 approaches behave against MICE, so we summarize here the key points:
- 941 • first and most importantly, there is not one single type of our models that perform better
 942 than the other. Both Generative Forests and Ensembles of Generative Trees can be useful
 943 for the task of missing data imputation. For example, *analcatdata_supreme* gives a clear
 944 advantage to Generative Forests: they even beat MICE with random forests having 4 000
 945 trees (that is hundreds of times more than our models) on both categorical variables (*perr*)
 946 and numerical variables (*rmse*). However, on *circgauss*, *student_performance_por* and
 947 *sigma_cabs*, Ensembles of Generative Trees obtain the best results. On *sigma_cabs*,
 948 they can perform on par with MICE if the number of tree is limited (which represents 100+
 949 times less trees than MICE's models);
- 950 • as is already noted in the main file (MF) and visible on several plots (see *e.g.* *sigma_cabs*),
 951 overfitting can happen with our models (both Generative Forests and Ensembles of Gen-
 952 erative Trees), which supports the idea that a pruning mechanism or a mechanism to stop
 953 training would be a strong plus to training such models. Interestingly, in several cases where

- 954 overfitting seems to happen, increasing the number of splits can reduce the phenomenon, at
 955 fixed number of trees: see for example `iris` ($50 \rightarrow 100$ iterations for EOGTs), `gridgauss`,
 956 `randgauss` and `kcc1` (GFs);
- 957 • some domains show that our models seem to be better at estimating continuous (re-
 958 gression) rather than categorical (prediction) variables: see for example `abalone`,
 959 `student_performance_mat`, `student_performance_por` and our biggest domain,
 960 `open_policing_hartford`. Note from that last domain that a larger number of itera-
 961 tions or models with a larger number of trees may be required for the best results, in
 962 particular for Ensembles of Generative Trees;
- 963 • small models may be enough to get excellent results on real world domains whose size
 964 would, at first glance, seem to require much bigger ones: on `compas`, we can compete (rmse)
 965 or beat (perr) `MICE` whose models contain 7 000 trees with just 20 stumps;
- 966 • it is important to remember that our technique does not just offer the capability to do missing
 967 data imputation: our models can also generate data and compute the full density conditional
 968 to observed values, which is not the case of `MICE`'s models, crafted with the sole purpose of
 969 solving the task of missing data imputation. Our objective was not to beat `MICE`, but rather
 970 show that Generative Forests and Ensembles of Generative Trees can also be useful for this
 971 task.

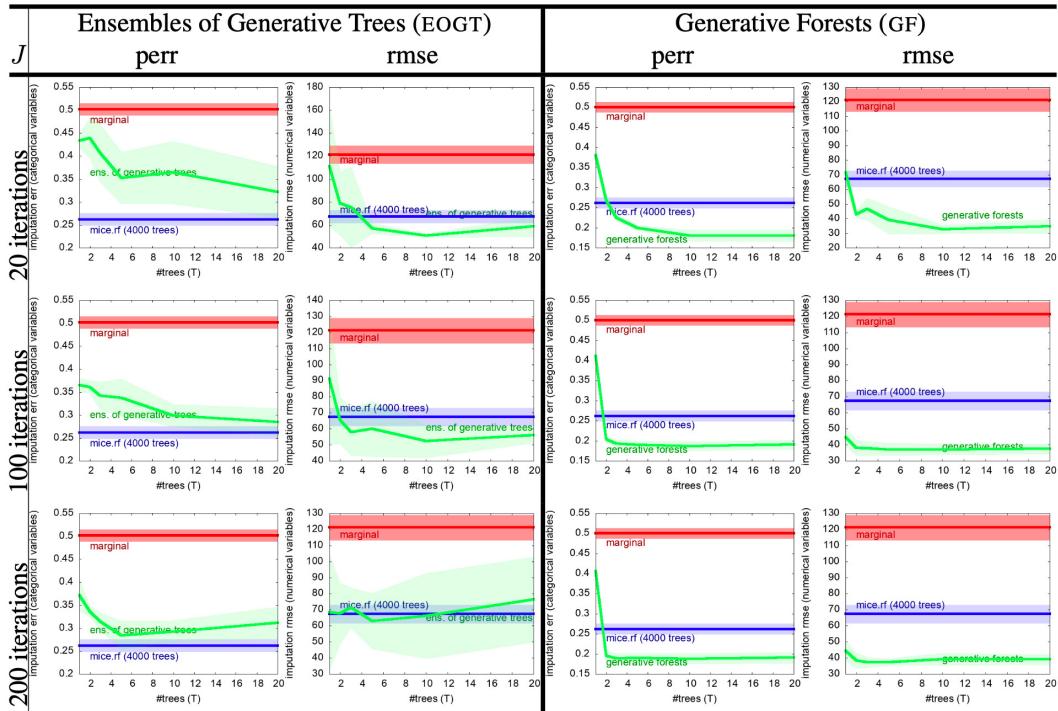


Table A13: Missing data imputation: results on `analcatdata_supreme` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs `mice` using RF (random forests) with 100 trees each. Each row corresponds to a different total number of iterations J in GF.BOOT. Since the domain contains both numerical and categorical variables, de compute separately the errors on categorical variables (using the error probability, denoted perr) and the error on numerical variables (using the root mean square error, rmse). In each plot, we display both the corresponding results of GF.BOOT, the results of `mice` (indicating also the total number of trees used to impute the dataset) and the result of a fast imputation baseline, the marginal computed from the training sample. In each plot, the x axis displays the number of trees T in GF.BOOT and each curve is displayed with its average \pm standard deviation in shaded color. The result for $T = 1$ equivalently indicates the performance of a single generative tree (GT) with J splits [29].

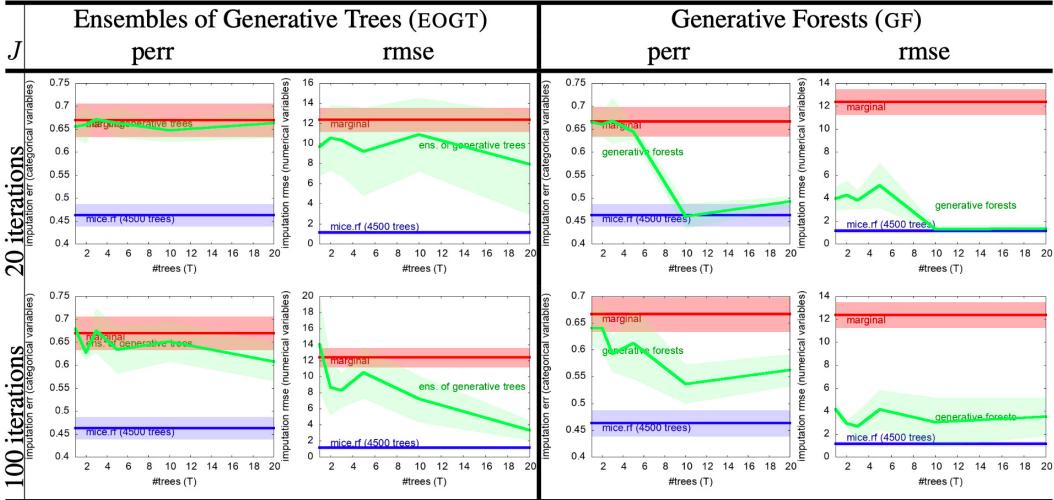


Table A14: Missing data imputation: results on abalone with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs mice using RF (random forests). Conventions follow Table A13.

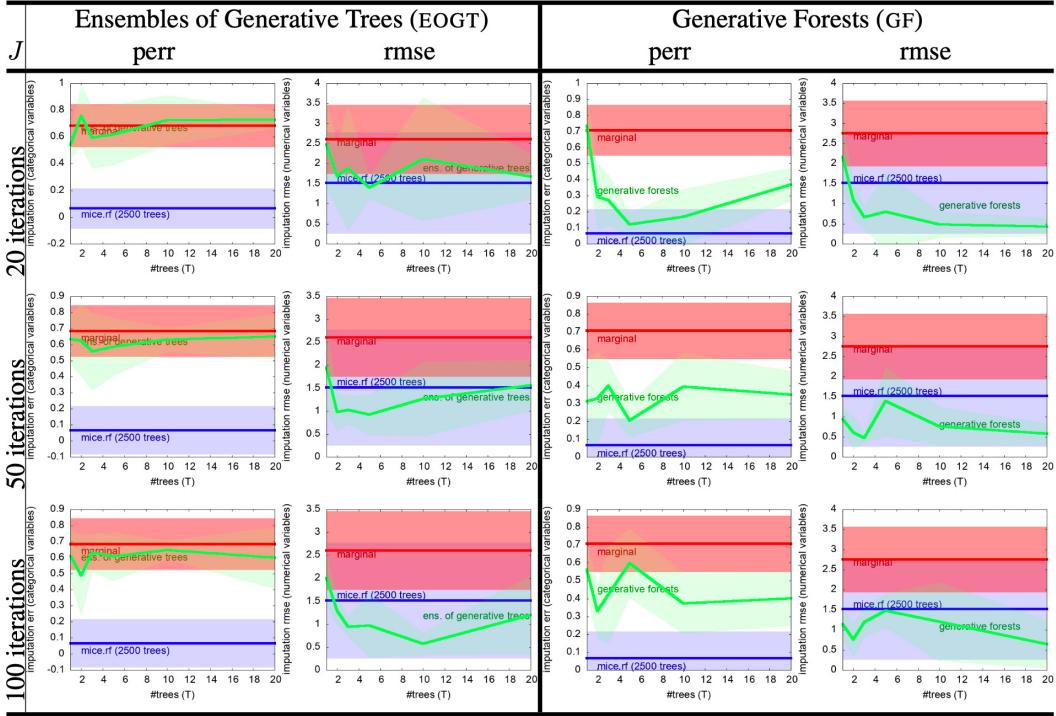


Table A15: Missing data imputation: results on `iris` with 5% missing features (MCAR), for GF BOOST learning GFs and EOGTs, vs `mice` using RF (random forests). Conventions follow Table A13.

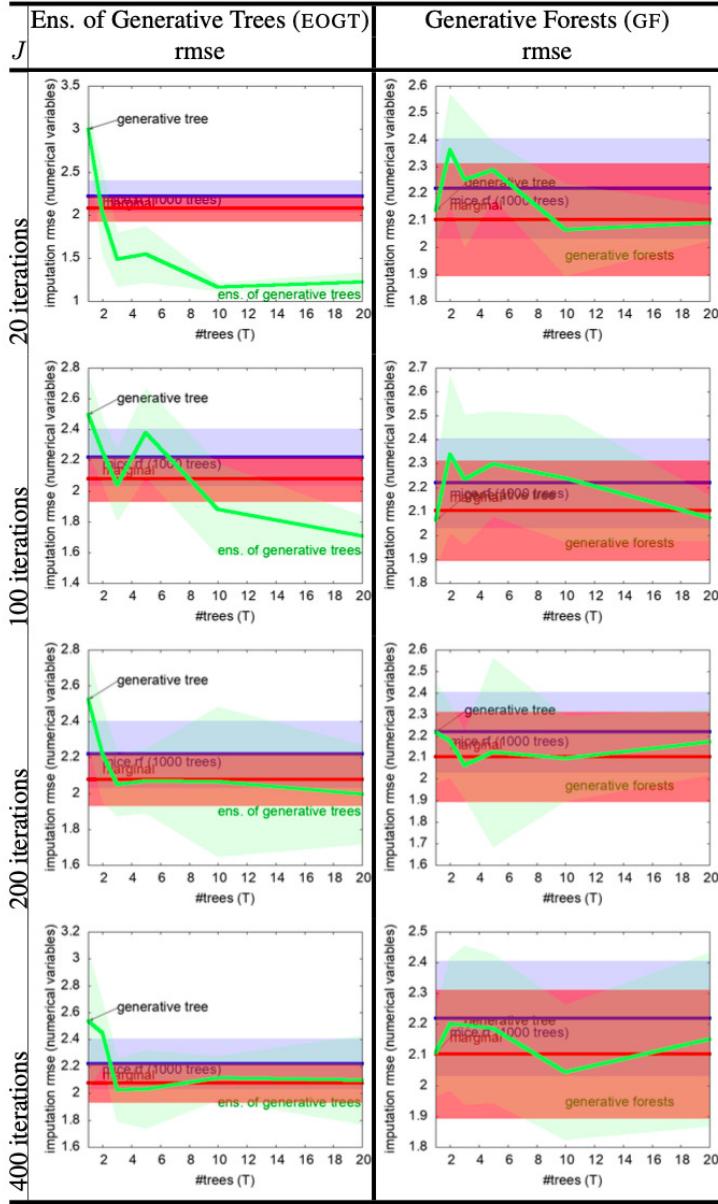


Table A16: Missing data imputation: results on `ringgauss` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs `mice` using RF (random forests). Since there are no nominal attributes in the domain, we have removed column `perr`. Conventions otherwise follow Table A13.

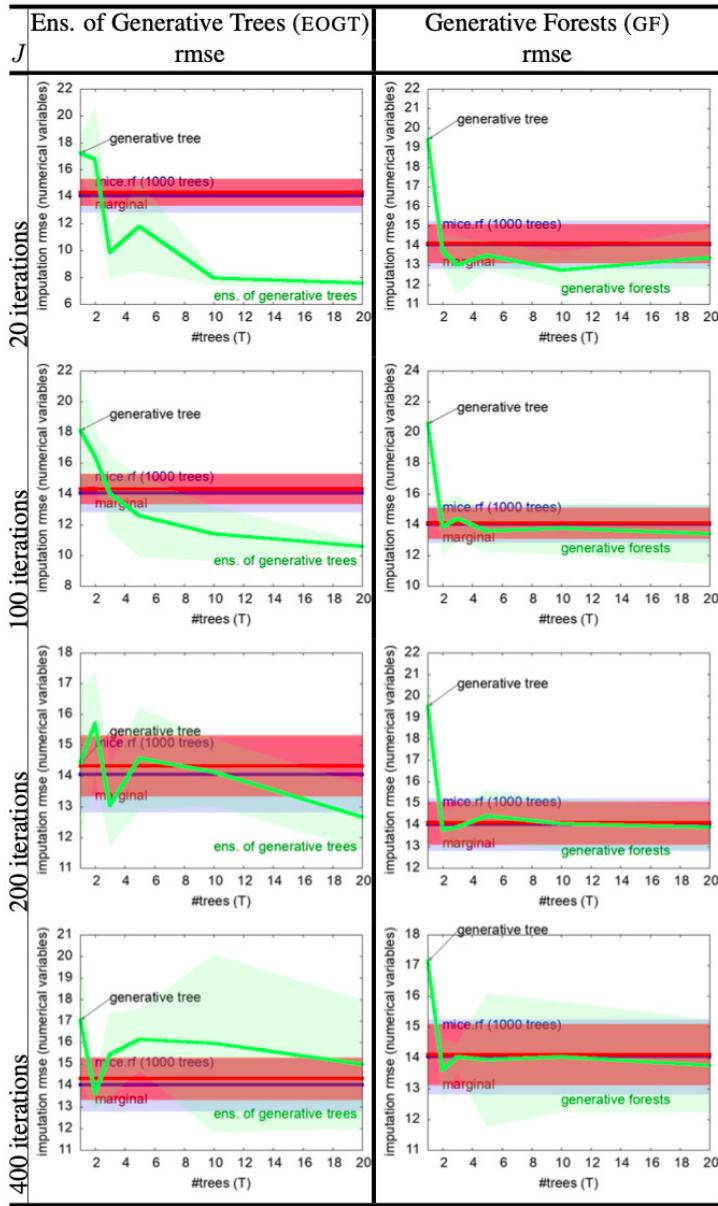


Table A17: Missing data imputation: results on `circgauss` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs `mice` using RF (random forests). Since there are no nominal attributes in the domain, we have removed column `perr`. Conventions otherwise follow Table A13.

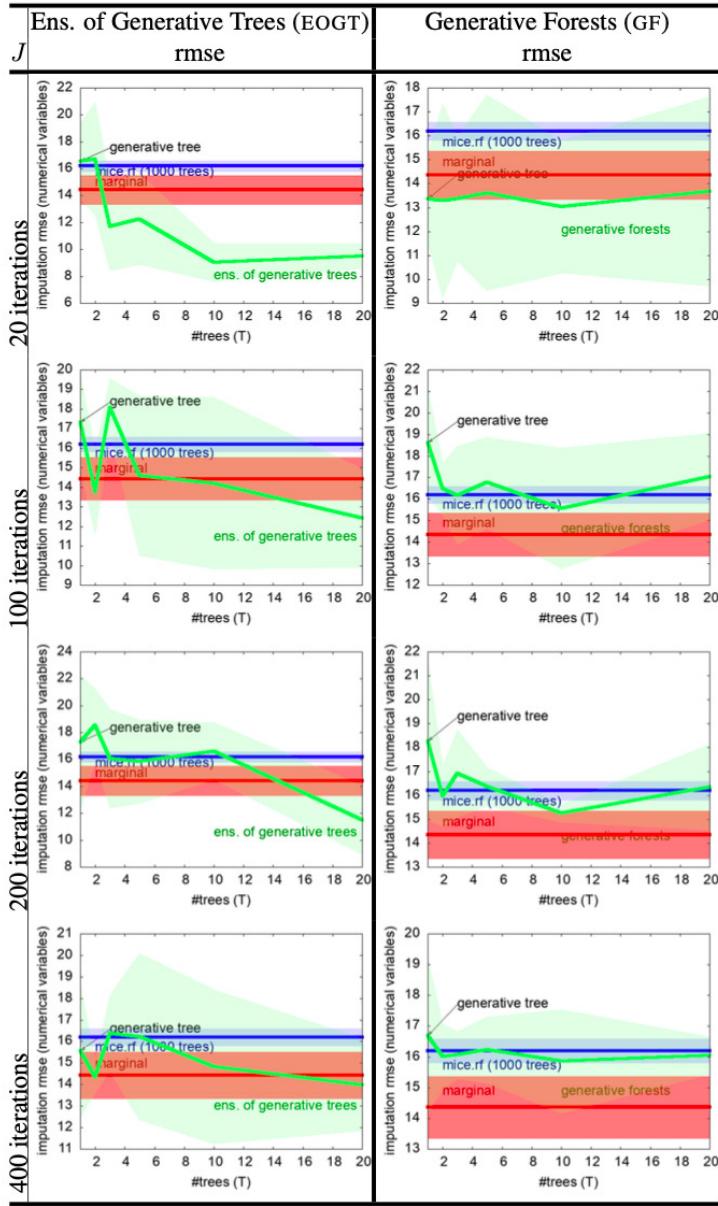


Table A18: Missing data imputation: results on gridgauss with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs mice using RF (random forests). Since there are no nominal attributes in the domain, we have removed column perr. Conventions otherwise follow Table A13.

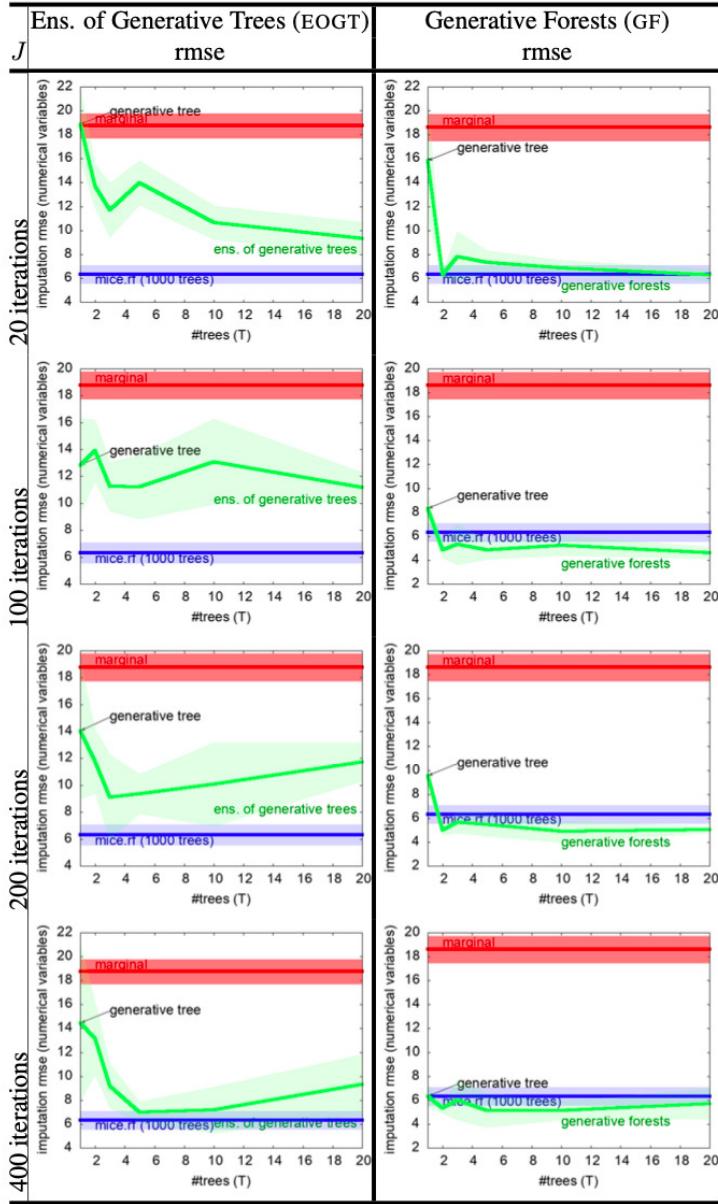


Table A19: Missing data imputation: results on `randgauss` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs `mice` using RF (random forests). Since there are no nominal attributes in the domain, we have removed column `perr`. Conventions otherwise follow Table A13.

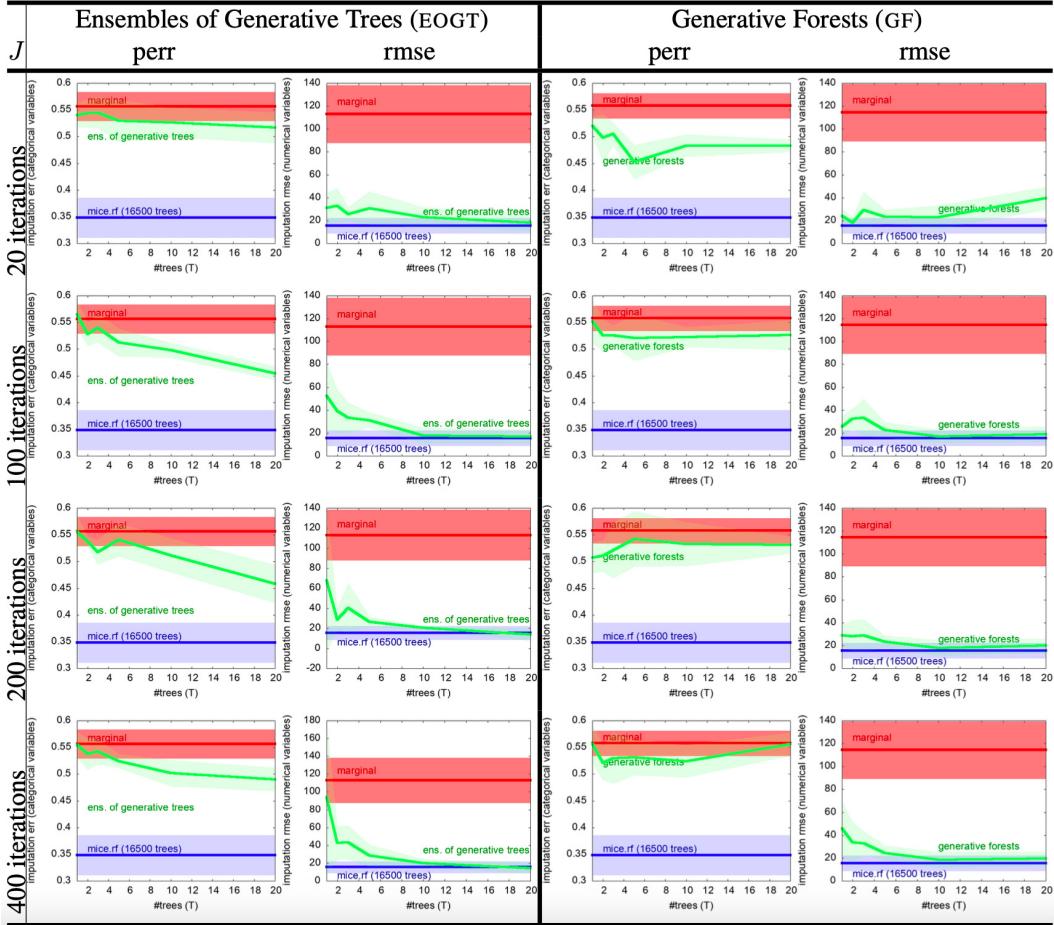


Table A20: Missing data imputation: results on `student_performance_mat` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs `mice` using RF (random forests). Conventions follow Table A13.

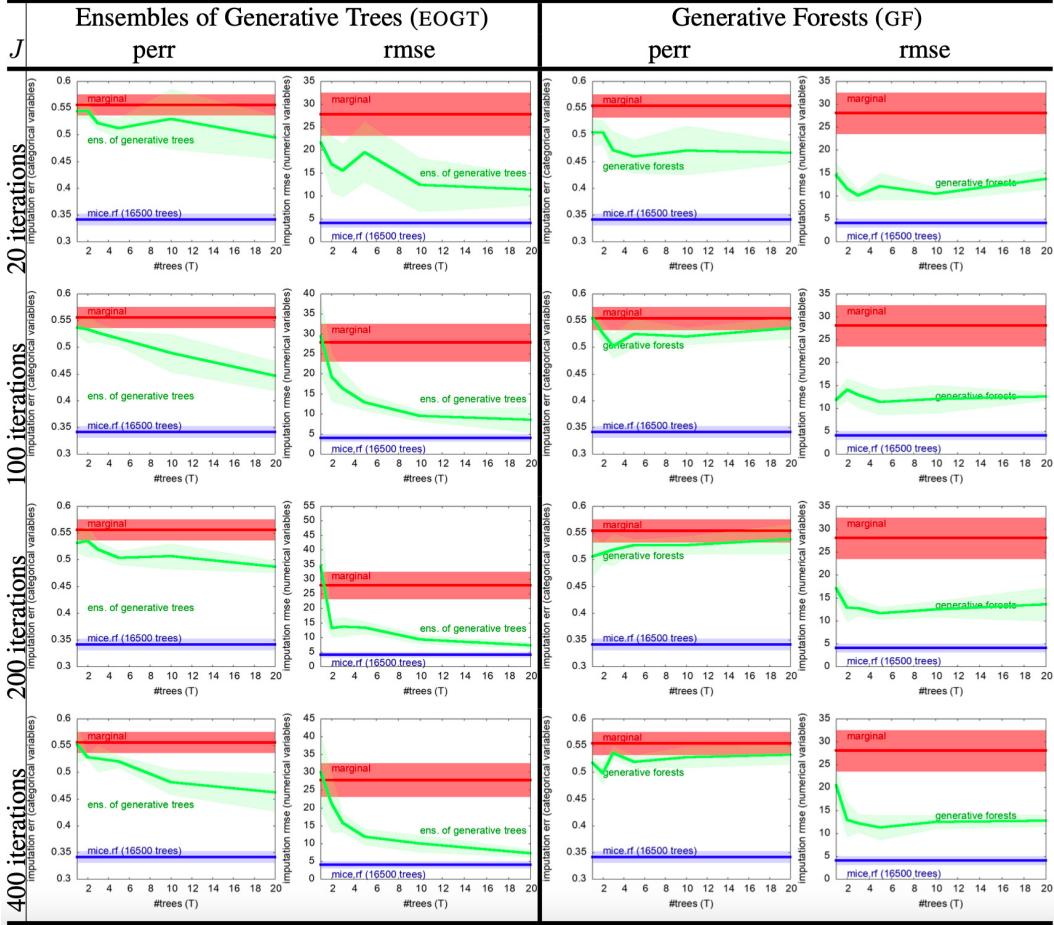


Table A21: Missing data imputation: results on `student_performance_por` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs mice using RF (random forests). Conventions follow Table A13.

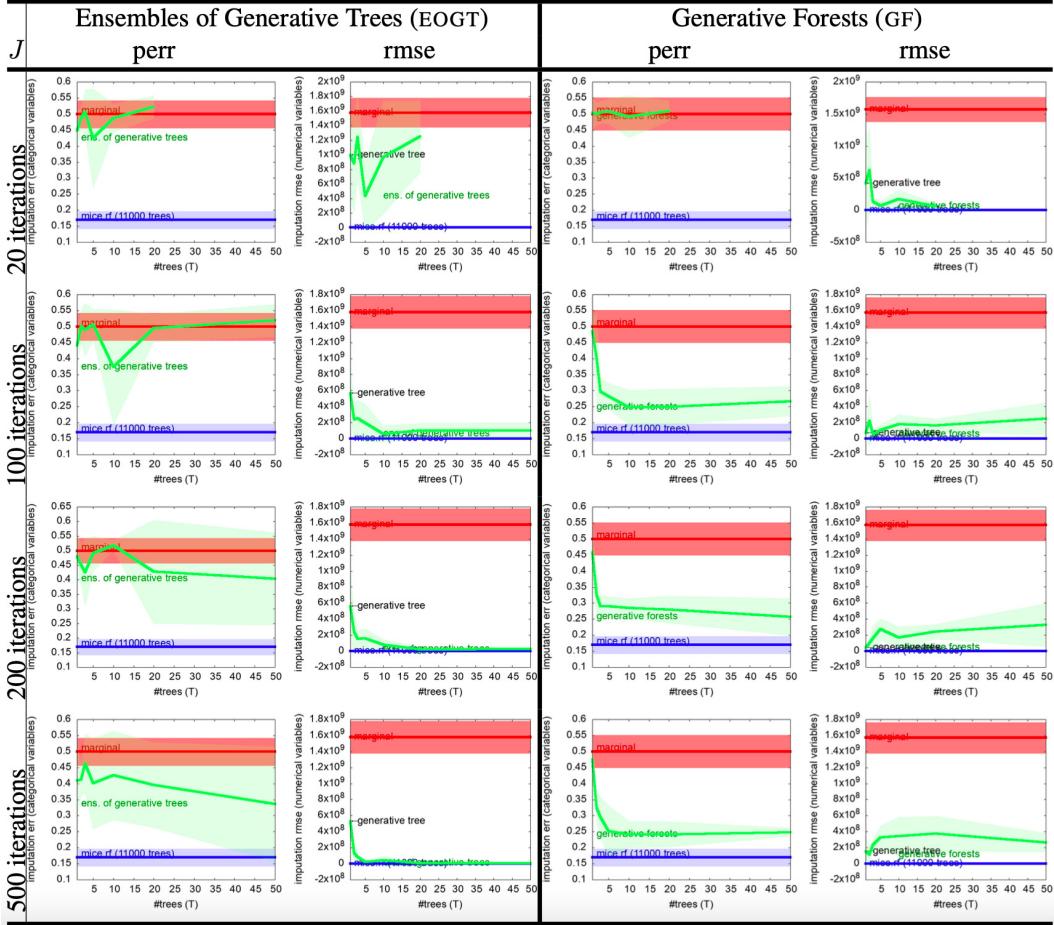


Table A22: Missing data imputation: results on kc1 with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs mice using RF (random forests). Conventions follow Table A13.

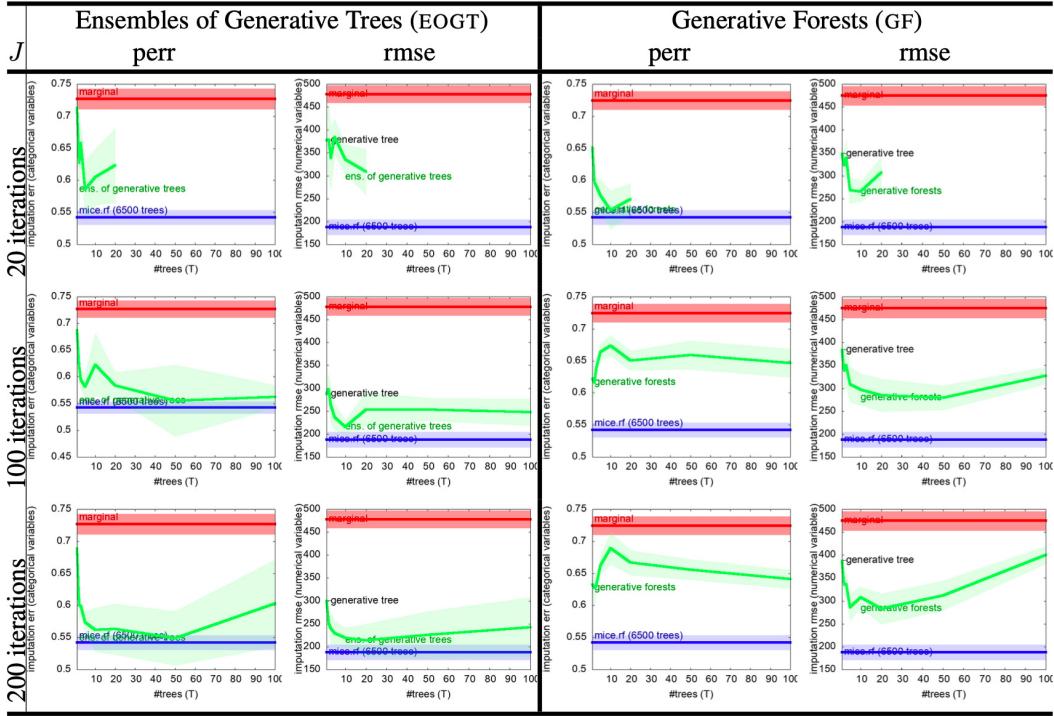


Table A23: Missing data imputation: results on `sigma_cabs` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs `mice` using RF (random forests). Conventions follow Table A13.

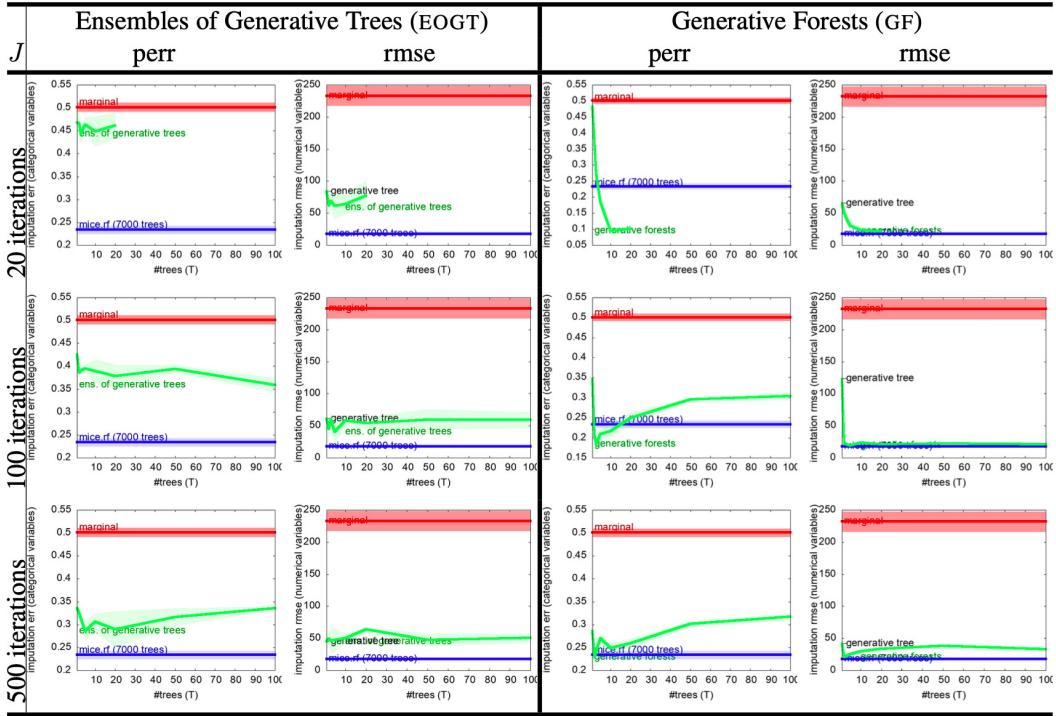


Table A24: Missing data imputation: results on `compas` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs `mice` using RF (random forests). Conventions follow Table A13.

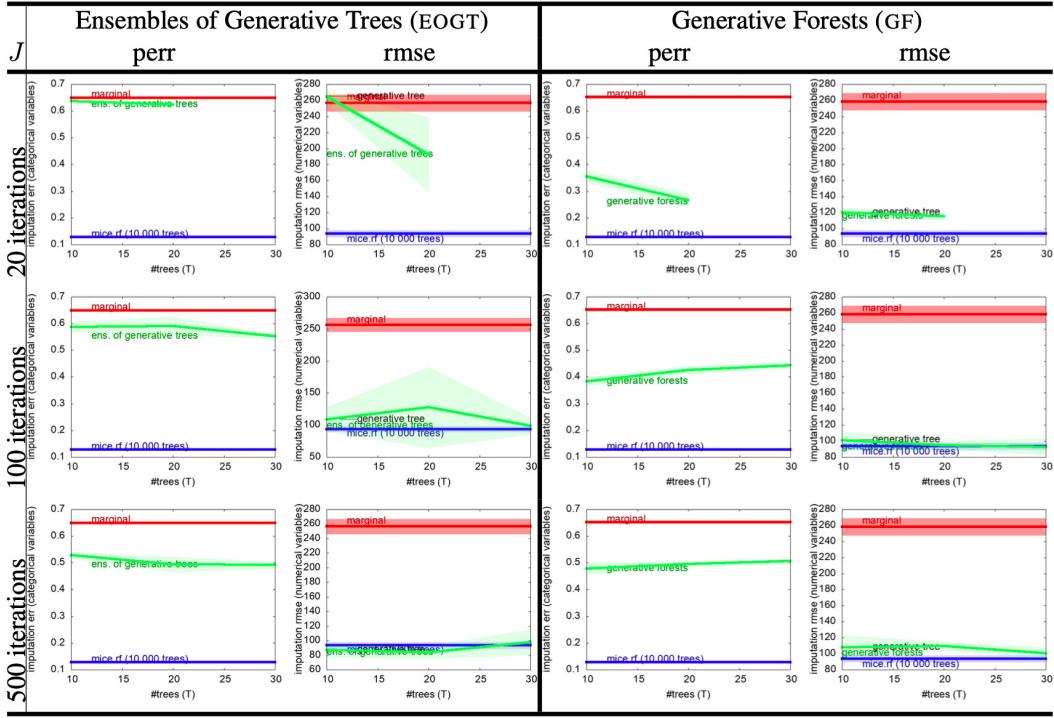


Table A25: Missing data imputation: results on `open_policing_hartford` with 5% missing features (MCAR), for GF.BOOT learning GFs and EOGTs, vs mice using RF (random forests). Conventions follow Table A13.

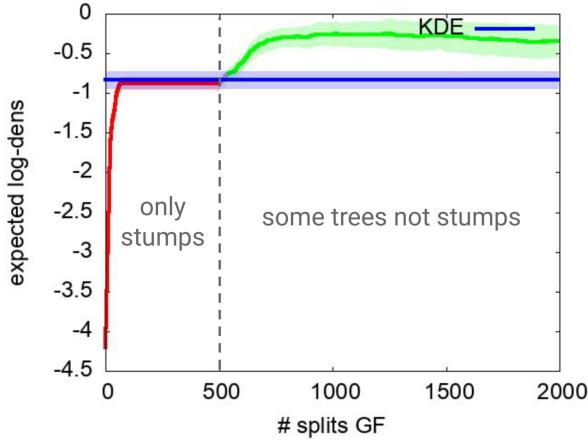


Table A26: DENSITY: zoom on the gridGauss domain: the GF is grown by GF.BOOST in a way that first learns a complete set of stumps (there are a total of $T = 500$ trees in the final model) and then grows some of the trees. The plateau for $j \leq 500$ leaves displays that we quickly reach a limit in terms of maximal performances of stumps, this being probably due to the fact that the dimension of the domain is small (2). We can however remark that with just stumps, we still attain close performances to KDE (note that we plot the expected log-densities; see text for details).

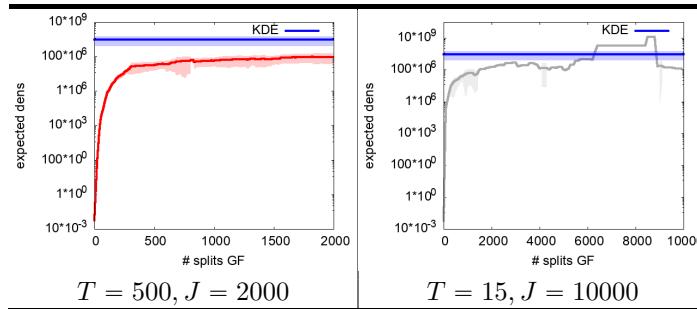


Table A27: DENSITY: zoom on the abalone domain: with $T = 500$ trees, total of $J = 2000$ as in Tables A28 and A29, we are clearly beaten by KDE (left). A much smaller of bigger trees ($T = 15, J = 10000$) is enough to become competitive and beat KDE, albeit not statistically significantly (right). The right picture also displays the interest into getting pruning or early stopping algorithms to prevent eventual overfitting.

972 V.V.3.7 Full comparisons with KDE on density estimation

973 For each domain we either use the expected density, or the expected negative log-density, as the metric
 974 to be maximized. We consider the expected density to cope with the eventuality of zero (pointwise)
 975 density predictions. Tables A28 and A29 present the results against KDE (details in Table A28).
 976 Table A26 singles out a result on the domain gridGauss, which displays an interesting pattern. The
 977 leaf chosen for a split in GF.BOOST is always the heaviest leaf; hence, as long as the iteration number
 978 $j \leq T$ (the maximum number of trees, which is 500 in this experiment), GF.BOOST grows stumps.
 979 In the case of gridGauss, like in a few other domains, there is a clear plateau in performances for
 980 $j \leq 500$, which does not hold anymore as $j > 500$. This, we believe, characterizes the fact that on
 981 these domains, which all share the property that their dimension is small, training stumps attains a
 982 limit in terms of performances. Table A27 shows that by changing basic parameters – here, decreasing
 983 the number of trees, increasing the total number of iterations –, one can obtain substantially better
 984 results.

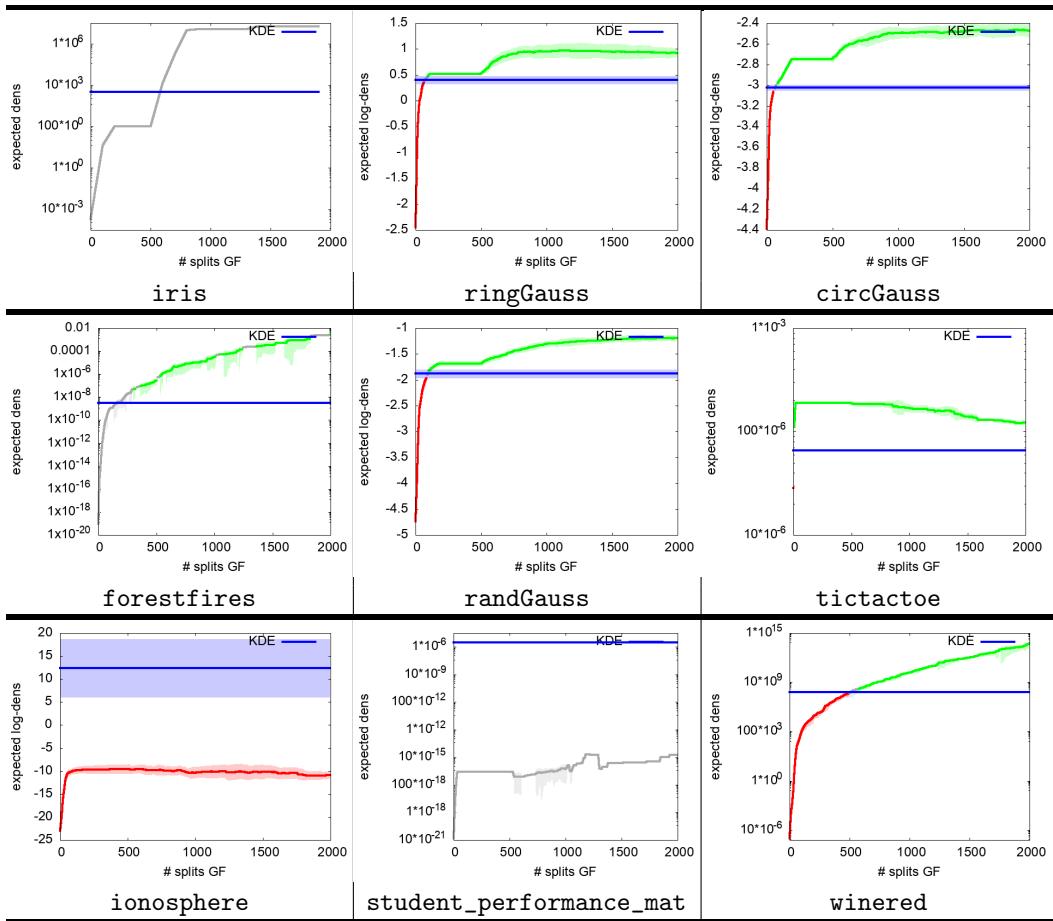


Table A28: Comparison of results for Kernel Density Estimation (KDE) and Generative Forests (us, GF) with parameters $T = 500$ trees, total of $J = 2000$ splits in trees, on a first batch of domains, put in increasing size according to Table A1. For each domain, we compute for GF at regular intervals the average \pm standard deviation of either (i) density values or (ii) negative log density values for the test fold (sometimes, KDE gives 0 estimated density which prevents the computation of (ii)). Warning: some y scales are logscale. We then compute the baseline of KDE and display its average \pm standard deviation in blue on each picture. We perform, at each applicable iteration of GF.BOOST (*i.e.* each iteration for many domains, which can be seen from the plots), a Student paired t -test over the five folds to compare the results with KDE. If we are statistically better ($p = 0.1$), our plot is displayed in green; if KDE is better, our plot is displayed in red; otherwise it is displayed in gray. When standard deviations are not displayed, it means average \pm standard deviation exceeds the plot's y scale.

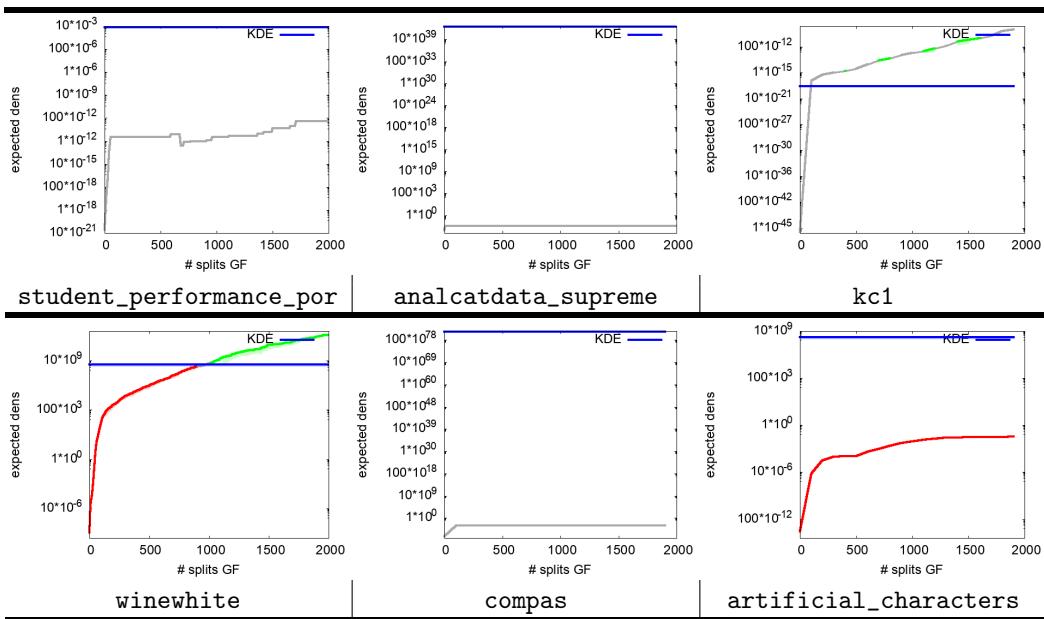


Table A29: Comparison of results for Kernel Density Estimation and Generative Forests (us, GF) where the number of trees and number of iterations are ($T = 500$ trees, total of $J = 2000$ splits in trees), on a second batch of domains. Conventions follow Table A28.

985 **NeurIPS Paper Checklist**

986 **1. Claims**

987 Question: Do the main claims made in the abstract and introduction accurately reflect the
988 paper's contributions and scope?

989 Answer: [Yes]

990 Justification: Claims are supported by both the theoretical results presented in the paper and
991 extensive experimental results, summarized in the main paper but otherwise presented in
992 extenso in the Supplement.

993 Guidelines:

- 994 • The answer NA means that the abstract and introduction do not include the claims
995 made in the paper.
- 996 • The abstract and/or introduction should clearly state the claims made, including the
997 contributions made in the paper and important assumptions and limitations. A No or
998 NA answer to this question will not be perceived well by the reviewers.
- 999 • The claims made should match theoretical and experimental results, and reflect how
1000 much the results can be expected to generalize to other settings.
- 1001 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1002 are not attained by the paper.

1003 **2. Limitations**

1004 Question: Does the paper discuss the limitations of the work performed by the authors?

1005 Answer: [Yes]

1006 Justification: Limitations are discussed, in particular in the experimental section vs the
1007 SOTA (see e.g. experiment Lifelike).

1008 Guidelines:

- 1009 • The answer NA means that the paper has no limitation while the answer No means that
1010 the paper has limitations, but those are not discussed in the paper.
- 1011 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1012 • The paper should point out any strong assumptions and how robust the results are to
1013 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1014 model well-specification, asymptotic approximations only holding locally). The authors
1015 should reflect on how these assumptions might be violated in practice and what the
1016 implications would be.
- 1017 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1018 only tested on a few datasets or with a few runs. In general, empirical results often
1019 depend on implicit assumptions, which should be articulated.
- 1020 • The authors should reflect on the factors that influence the performance of the approach.
1021 For example, a facial recognition algorithm may perform poorly when image resolution
1022 is low or images are taken in low lighting. Or a speech-to-text system might not be
1023 used reliably to provide closed captions for online lectures because it fails to handle
1024 technical jargon.
- 1025 • The authors should discuss the computational efficiency of the proposed algorithms
1026 and how they scale with dataset size.
- 1027 • If applicable, the authors should discuss possible limitations of their approach to
1028 address problems of privacy and fairness.
- 1029 • While the authors might fear that complete honesty about limitations might be used by
1030 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1031 limitations that aren't acknowledged in the paper. The authors should use their best
1032 judgment and recognize that individual actions in favor of transparency play an impor-
1033 tant role in developing norms that preserve the integrity of the community. Reviewers
1034 will be specifically instructed to not penalize honesty concerning limitations.

1035 **3. Theory Assumptions and Proofs**

1036 Question: For each theoretical result, does the paper provide the full set of assumptions and
1037 a complete (and correct) proof?

1038 Answer: [Yes]

1039 Justification: Assumptions are detailed, formal results given and full proofs are provided in
1040 the Appendix.

1041 Guidelines:

- 1042 • The answer NA means that the paper does not include theoretical results.
1043 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
1044 referenced.
1045 • All assumptions should be clearly stated or referenced in the statement of any theorems.
1046 • The proofs can either appear in the main paper or the supplemental material, but if
1047 they appear in the supplemental material, the authors are encouraged to provide a short
1048 proof sketch to provide intuition.
1049 • Inversely, any informal proof provided in the core of the paper should be complemented
1050 by formal proofs provided in appendix or supplemental material.
1051 • Theorems and Lemmas that the proof relies upon should be properly referenced.

1052 **4. Experimental Result Reproducibility**

1053 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
1054 perimental results of the paper to the extent that it affects the main claims and/or conclusions
1055 of the paper (regardless of whether the code and data are provided or not)?

1056 Answer: [Yes]

1057 Justification: (i) full code provided, (ii) extensive README.txt, (iii) coding choices dis-
1058 cussed in Appendix, (iv) all scripts provided for easy running of each experiments, (v)
1059 special visualization routines (e.g. heatmaps) also included in the code, (vi) example
1060 domains (public and simulated) provided for quick assessment.

1061 Guidelines:

- 1062 • The answer NA means that the paper does not include experiments.
1063 • If the paper includes experiments, a No answer to this question will not be perceived
1064 well by the reviewers: Making the paper reproducible is important, regardless of
1065 whether the code and data are provided or not.
1066 • If the contribution is a dataset and/or model, the authors should describe the steps taken
1067 to make their results reproducible or verifiable.
1068 • Depending on the contribution, reproducibility can be accomplished in various ways.
1069 For example, if the contribution is a novel architecture, describing the architecture fully
1070 might suffice, or if the contribution is a specific model and empirical evaluation, it may
1071 be necessary to either make it possible for others to replicate the model with the same
1072 dataset, or provide access to the model. In general, releasing code and data is often
1073 one good way to accomplish this, but reproducibility can also be provided via detailed
1074 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1075 of a large language model), releasing of a model checkpoint, or other means that are
1076 appropriate to the research performed.
1077 • While NeurIPS does not require releasing code, the conference does require all submis-
1078 sions to provide some reasonable avenue for reproducibility, which may depend on the
1079 nature of the contribution. For example
1080 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1081 to reproduce that algorithm.
1082 (b) If the contribution is primarily a new model architecture, the paper should describe
1083 the architecture clearly and fully.
1084 (c) If the contribution is a new model (e.g., a large language model), then there should
1085 either be a way to access this model for reproducing the results or a way to reproduce
1086 the model (e.g., with an open-source dataset or instructions for how to construct
1087 the dataset).
1088 (d) We recognize that reproducibility may be tricky in some cases, in which case
1089 authors are welcome to describe the particular way they provide for reproducibility.
1090 In the case of closed-source models, it may be that access to the model is limited in
1091 some way (e.g., to registered users), but it should be possible for other researchers
1092 to have some path to reproducing or verifying the results.

1093 **5. Open access to data and code**

1094 Question: Does the paper provide open access to the data and code, with sufficient instructions
1095 to faithfully reproduce the main experimental results, as described in supplemental
1096 material?

1097 Answer: [Yes]

1098 Justification: **(i)** full code provided, **(ii)** extensive README.txt, **(iii)** coding choices dis-
1099 cussed in Appendix, **(iv)** all scripts provided for easy running of each experiments, **(v)**
1100 special visualization routines (e.g. heatmaps) also included in the code, **(vi)** example
1101 domains (public and simulated) provided for quick assessment.

1102 Guidelines:

- 1103 • The answer NA means that paper does not include experiments requiring code.
- 1104 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1105 • While we encourage the release of code and data, we understand that this might not be
1106 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
1107 including code, unless this is central to the contribution (e.g., for a new open-source
1108 benchmark).
- 1109 • The instructions should contain the exact command and environment needed to run to
1110 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1111 • The authors should provide instructions on data access and preparation, including how
1112 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1113 • The authors should provide scripts to reproduce all experimental results for the new
1114 proposed method and baselines. If only a subset of experiments are reproducible, they
1115 should state which ones are omitted from the script and why.
- 1116 • At submission time, to preserve anonymity, the authors should release anonymized
1117 versions (if applicable).
- 1118 • Providing as much information as possible in supplemental material (appended to the
1119 paper) is recommended, but including URLs to data and code is permitted.

1122 **6. Experimental Setting/Details**

1123 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
1124 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
1125 results?

1126 Answer: [Yes]

1127 Justification: Details are provided in main text and Appendix.

1128 Guidelines:

- 1129 • The answer NA means that the paper does not include experiments.
- 1130 • The experimental setting should be presented in the core of the paper to a level of detail
1131 that is necessary to appreciate the results and make sense of them.
- 1132 • The full details can be provided either with the code, in appendix, or as supplemental
1133 material.

1134 **7. Experiment Statistical Significance**

1135 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1136 information about the statistical significance of the experiments?

1137 Answer: [Yes]

1138 Justification: standard devs and inferential statistics test done (details in main file + ap-
1139 pendix).

1140 Guidelines:

- 1141 • The answer NA means that the paper does not include experiments.
- 1142 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
1143 dence intervals, or statistical significance tests, at least for the experiments that support
1144 the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA] .

Justification: Our code runs on any standard computer. Details of codes, computer and SOTA used (version, etc.) in appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research of the paper follows the code of ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: See conclusion.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No release of data or models.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Appropriate credit / references are provided. Licenses listed in Appendix.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 1247 • If assets are released, the license, copyright information, and terms of use in the
1248 package should be provided. For popular datasets, paperswithcode.com/datasets
1249 has curated licenses for some datasets. Their licensing guide can help determine the
1250 license of a dataset.
1251 • For existing datasets that are re-packaged, both the original license and the license of
1252 the derived asset (if it has changed) should be provided.
1253 • If this information is not available online, the authors are encouraged to reach out to
1254 the asset's creators.

1255 **13. New Assets**

1256 Question: Are new assets introduced in the paper well documented and is the documentation
1257 provided alongside the assets?

1258 Answer: [NA]

1259 Justification: No new assets provided.

1260 Guidelines:

- 1261 • The answer NA means that the paper does not release new assets.
1262 • Researchers should communicate the details of the dataset/code/model as part of their
1263 submissions via structured templates. This includes details about training, license,
1264 limitations, etc.
1265 • The paper should discuss whether and how consent was obtained from people whose
1266 asset is used.
1267 • At submission time, remember to anonymize your assets (if applicable). You can either
1268 create an anonymized URL or include an anonymized zip file.

1269 **14. Crowdsourcing and Research with Human Subjects**

1270 Question: For crowdsourcing experiments and research with human subjects, does the paper
1271 include the full text of instructions given to participants and screenshots, if applicable, as
1272 well as details about compensation (if any)?

1273 Answer: [NA]

1274 Justification: No crowdsourcing or research with human subjects.

1275 Guidelines:

- 1276 • The answer NA means that the paper does not involve crowdsourcing nor research with
1277 human subjects.
1278 • Including this information in the supplemental material is fine, but if the main contribu-
1279 tion of the paper involves human subjects, then as much detail as possible should be
1280 included in the main paper.
1281 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1282 or other labor should be paid at least the minimum wage in the country of the data
1283 collector.

1284 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
1285 Subjects**

1286 Question: Does the paper describe potential risks incurred by study participants, whether
1287 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1288 approvals (or an equivalent approval/review based on the requirements of your country or
1289 institution) were obtained?

1290 Answer: [NA]

1291 Justification: No research with human subjects.

1292 Guidelines:

- 1293 • The answer NA means that the paper does not involve crowdsourcing nor research with
1294 human subjects.
1295 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1296 may be required for any human subjects research. If you obtained IRB approval, you
1297 should clearly state this in the paper.

- 1298
- We recognize that the procedures for this may vary significantly between institutions
- 1299
- and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1300
- guidelines for their institution.
- 1301
- For initial submissions, do not include any information that would break anonymity (if
- 1302
- applicable), such as the institution conducting the review.