

1 We want to thank all reviewers for their time and effort. We address the raised points below.

- 2 – **R1** *hard to learn large number of densities*: We used a threshold on the number of samples at each leaf, to decide
3 whether LearnSPN or a fully factorised density is used. This drastically reduces the computational effort to learn
4 the leaves. In general, GeDTs/GeFs allow for much freedom, as any density estimator can be used at the leaves.
- 5 – **R1** *additional computational overhead*: The asymptotic costs of learning and inference are comparable to standard
6 RFs (see supp. material). In general, the overhead depends on the particular choice of leaf estimators.
- 7 – **R1** *when to use SPNs or fully factorised leaves*: A key consideration is overfitting, and using factorised leaves
8 induces a conditional (on the DT decisions) independence assumption which is more robust for few samples.
- 9 – **R1** *prediction goes to 0 as X_2 approaches $E(X_2|X_1)$ instead of 1*: This is indeed a typo: Prediction will be (w.h.p.)
10 $Y = 0$ (the argumentation remains valid as presented in the paper). Thanks for the catch!
- 11 – **R1** *Fig.2 is slightly puzzling*: We see the point, and will revise the figure.
- 12 – **R3** *compare against the most related approaches (e.g., DT+KDE, DT+NBC, C Nets)*: Note that GeDTs/GeFs
13 leaves can be equipped with any density model. Therefore, when we equip them with KDEs or NBCs, they simply
14 subsume DT+KDE and DT+NBC in terms of classification. The crucial difference in our paper is, that these
15 previous approaches did not follow through with the interpretation of the **whole** DT as a joint distribution and its
16 advantages (missing data, outlier detection). We will include a comparison to C Nets.
- 17 – **R5** *practical interest beyond handling missing values, properties of GeDT/GeRF*: We focused on classification under
18 missing data since it is the most compelling advantage: DT/RF practitioners do not need to change the structure
19 learning algorithm, but can now treat missing inputs within the same model, and with guaranteed backwards
20 compatibility (provided $p(Y, \mathbf{X}) = p(Y)p(\mathbf{X})$ in all leaves). We will include more thorough experiments on outlier
21 detection in the revised version. Our main contribution is that we link two separate communities, PCs and DTs,
22 which will likely lead to ample cross-fertilisation of ideas.
- 23 – **R5** *RF hyper-parameters*: Our experiments only compare models with the same underlying structure, so variations
24 in the RF hyper-parameters would most likely only change the basis for comparison but not the overall results.
25 As our goal is not to find the most accurate RF/GeF, we use literature defaults and defer tuning to future work
26 (non-tuned results are already very good). That said, we did experiment with a few hyper-parameters. For instance,
27 pruning has not proved very useful (not uncommon in RFs), see right side of Figure below.
- 28 – **R5** *LearnSPN on the full dataset*: We experimented with a SOTA implementation of LearnSPN and results are
29 quite poor. For the sake of space we omitted it in the submission (see left side of Figure below), but will include it
30 given acceptance.
- 31 – **R5** *explanations on density model part*: GeDTs/GeFs can be used with any density estimator, which is subject to
32 exploration in the future (we use SPNs as cited). Also, we will explain fully factorised leaves in the paper. $P(X_i)$
33 is a marginal (multinomial/normal) trained on data matching the leaf.
- 34 – **R5** *Friedman method*: We will clarify differences in the paper; results do not automatically apply, but it may be
35 possible to extend them. Note that Friedman’s method ignores the actual values of explanatory variables, while
36 fully factorised leaves do compute $p(\mathbf{x})$ to handle missing values.
- 37 – **R5** *MissForest, MERCS, Isolation Forest*: We agree that claiming GeFs are the new SOTA for treating missing
38 values would require further experimentation, but we do not make such a claim. In our experiments, we focus on
39 ‘built-in’ methods as they are closer to ours in nature, and also compared to KNN because it is the most common
40 imputation method and the closest to ours in computational cost. We will cite and discuss the works pointed out.

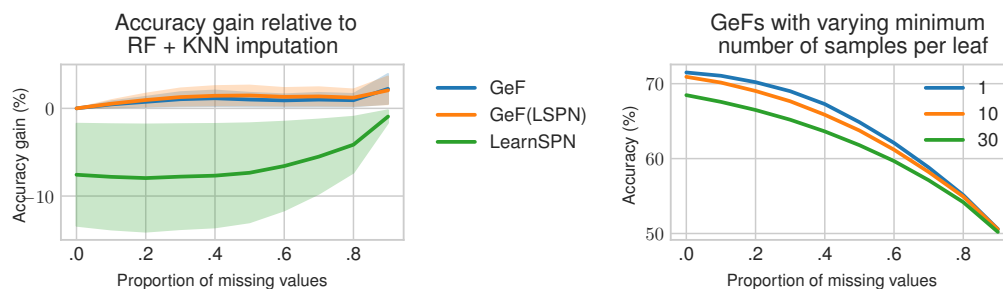


Figure 1: Average accuracy gain relative to RF+KNN imputation (left) and average accuracy on test data of GeF models with different minimum number of samples per leaf (right). Both averages are across the datasets considered.