# **Estimating Information-Theoretic Quantities with Uncertainty Forest**

Ronak Mehta<sup>1</sup>, Richard Guo<sup>1</sup>, Jesús Arroyo<sup>1</sup>, Mike Powell<sup>1</sup>, Hayden Helm<sup>1</sup>, Cencheng Shen<sup>1</sup>, and Joshua T. Vogelstein<sup>1,2\*</sup>

Abstract. Information-theoretic quantities, such as conditional entropy and mutual information, are critical data summaries for quantifying uncertainty. Existing estimators for these quantities either have strong theoretical guarantees or effective performance in high-dimensional data, but not both. We propose a decision forest method, Uncertainty Forests (UF), which combines quantile regression forests, honest sampling, and a finite sample correction. We prove UF provides consistent estimates for these information-theoretic quantities, including in multivariate settings. Empirically, UF reduces finite sample bias and variance in a range of both low- and high-dimensional simulated settings for estimating posterior probabilities, conditional entropies, and mutual information. In a real-world connectome application, UF quantifies the uncertainty about neuron type given various cellular features in the Drosophila larva mushroom body, a key challenge for modern neuroscience.

1 Introduction Uncertainty quantification is a fundamental desiderata of statistical inference and data science. In supervised learning settings it is common to quantify uncertainty with either conditional entropy or mutual information (MI). Suppose we are given a pair of random variables (X,Y), where X is d-dimensional vector-valued and Y is a categorical variable of interest. Conditional entropy H(Y|X) measures the uncertainty in Y on average given the value of X. On the other hand, mutual information quantifies the shared information between X and Y. Although both parameters are readily estimated when X and Y are low-dimensional and "nicely" distributed, an important problem arises in measuring these quantities from higher-dimensional data in a nonparametric fashion [1].

We address the problem of high-dimensionality by proposing a decision forest method for estimating conditional entropy under the framework that X is any d-dimensional random vector and Y is categorical. Because we restrict Y to be categorical, we can easily compute the maximum-likelihood estimate of H(Y) and hence use our estimator to compute mutual information.

We present Uncertainty Forests (UF) to estimate these information-theoretic quantities. UF combines quantile regression forests [2] with honest sampling [3], and introduces a finite sample correction to improve performance while preserving consistency. We prove that our estimation technique is consistent (meaning that it converges in probability to the true estimates of conditional entropy and mutual information, under mild conditions on the joint distribution of X and Y), and demonstrate via simulations that UF performs well in both low- and high-dimensional settings when estimating conditional distributions, conditional entropy, and mutual information.

Finally, we provide a real-world application of our estimator by measuring information about the various properties of *Drosophila* neurons that are contained in neuron cell types. Scientific prior knowledge poses various relationships between particular features and the cell type, which are supported by the mutual information and conditional mutual information estimates of UF.

# 2 Background

**2.1 Problem Formulation** Suppose we are given two random variables X and Y with support sets  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. Let x,y denote specific values that the random variables take on and p(x),p(y) be the probabilities of X=x and Y=y. The unconditioned Shannon entropy of Y is  $H(Y)=-\sum_{y\in\mathcal{Y}}p(y)\log p(y)$ . Analogously, conditional entropy is  $H(Y|X)=\sum_{x\in\mathcal{X}}p(x)H(Y|X=x)=-\sum_{x\in\mathcal{X}}p(x)\sum_{y\in\mathcal{Y}}p(y|x)\log p(y|x)$ , where p(y|x) is the conditional probability that Y=y given X=x and H(Y|X=x) is the entropy of Y conditioned on X equaling particular value x. This quantity represents the average uncertainty in Y having observed X. In the case of a continuous random variable, the sum over the corresponding support is replaced with an integral, and probability mass functions are replaced by densities. For the remainder of this work, we assume that  $\mathcal{X}=\mathbb{R}^d$  and  $\mathcal{Y}=[K]=\{1,...,K\}$  for positive integers d and K. Mutual information, I(X;Y), can be computed

<sup>\*</sup>Corresponding author: jovo@jhu.edu; 1Johns Hopkins University, 2Progressive Learning

from conditional entropy I(X;Y) = H(Y) - H(Y|X) = H(X) - H(X|Y). Mutual information has many appealing properties, such as symmetry, and is widely used in data science applications [4].

**2.2 Related Work** A common approach to estimating mutual information relies on the 3-H principle [4], I(X;Y) = H(Y) + H(X) - H(X,Y), where  $H(X,Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log p(x,y)$  is the Shannon entropy of the pair (X,Y). Procedures typically estimate unconditional entropy for the three separate random variables. Examples of these include kernel density estimates and ensembles of k-NN estimators [5–8]. One method in particular, the KSG estimator, popular because of its excellent empirical performance, improves k-NN estimates via heuristics [9]. Other approaches include binning [10] and von Mises estimators [11]. Unfortunately, many modern datasets contain mixtures of continuous inputs X and discrete outputs Y. In these cases, few of the above methods work well, as while individual entropies H(X) and H(Y) are well-defined, H(Y,X) is either ill-defined or not easily estimated, thus rendering the 3-H approach intractable [4].

A recent approach, referred to as Mixed KSG, modifies the KSG estimator to improve its performance in various settings, including mixed continuous and categorical inputs [4]. Computing both mutual information and conditional entropy becomes difficult in higher dimensional data. Numerical summations or integration become computationally intractable, and nonparametric methods (for example, k-nearest neighbor, kernel density estimates, binning, Edgeworth approximation, likelihood ratio estimators) typically do not scale well with increasing dimensions [1, 12]. While a recently proposed neural network approach, MINE [1], addresses high-dimensional data, the method does not return an estimate of the posterior  $p(y \mid x)$ , which is useful in uncertainty quantification. Additionally, MINE requires that the network of choice be expressive enough so that the Donsker-Varadhan representation can be used to approximate the mutual information, an assumption that is difficult to inspect. Other theoretically analyzed deep approaches to mutual information estimation can make stringent assumptions on the distribution of learned weights [13]. Finally, in the interest of an "out-of-the-box" method, another common pitfall of deep approaches is sensitivity to hyperparameter choice.

Another popular ad-hoc approach relies on estimation of the posterior probability  $p(y \mid x)$ . Calibration methods, such as isotonic regression, map classification scores, whether they are scalar projections of the data in linear discriminant analysis or empirical posteriors from random forests, to "calibrated" posteriors using a held-out set. A thorough theoretical analysis of such methods has yet to be addressed [14].

**2.3 CART Random Forest** Classification and Regression Tree (CART) Random Forest is a robust, powerful algorithm that leverages ensembles of decision trees for classification and regression tasks [15]. In a study of over 100 classification problems, Férnandez-Delgado et al. [16] showed that random forests have the overall best performance when compared 178 other classifiers. Furthermore, random forests are highly scalable; efficient implementations can build a forest of 100 trees from 110 Gigabyte data (n = 10,000,000, d = 1000) in little more than an hour [17].

Random forest classifiers are instances of a bagging classifiers, in which the base classifiers are decision trees. A decision tree first learns a partition of feature space, then learns constant functions within each part to perform classification or regression. Precisely,  $\mathcal{L}$  is called a partition of feature space  $\mathcal{X}$ , if for every  $L, L' \in \mathcal{L}$  with  $L \neq L', L \cap L' = \varnothing$ , and  $\bigcup_{L \in \mathcal{L}} L = \mathcal{X}$ . This  $\mathcal{L}$  is learned on data  $\{X_i, Y_i\}_{i=1}^n$  by recursively splitting a randomly selected subsample along a single dimension of the input data based on an impurity measure [15], such as Gini impurity or entropy. The trees are grown until nodes reach a certain criterion (for example, a minimum number of samples). The bottom-most nodes are called "leaf nodes". Additionally, only a random subset of dimensions of X are considered for the split at each node.

Given a partition  $\mathcal{L}$ , let L(x) be the part of  $\mathcal{L}$  to which  $x \in \mathcal{X}$  belongs. Letting  $\mathbb{1}[\cdot]$  be the indicator function, a possible predictor function  $\hat{g}$  for classification is  $\hat{g}(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{i=1}^n \mathbb{1}[Y_i = y \text{ and } X_i \in L(x)]$ . This is the plurality vote among the training data in the leaf node of x. For a regression task, a corresponding predictor  $\hat{\mu}$  is  $\hat{\mu}(x) = \frac{\sum_{i=1}^n Y_i \cdot \mathbb{1}[X_i \in L(x)]}{\sum_{i=1}^n \mathbb{1}[X_i \in L(x)]}$ , which is the average Y value for the training

data in L(x). For a forest, B trees are learned on randomly subsampled points. These points used in tree construction are called the 'in-bag' samples for that tree, while those that are left out are called 'out-of-bag' samples.

**2.4** Honest Sampling, Balanced Trees, and Random-Split Trees In decision forest algorithms, because data are typically split to maximize purity within child nodes, the posteriors estimated in each cell tend to be biased toward certainty. Honesty, as shown in Breiman [18], helps in bounding the bias of tree-based estimates of posterior class probabilities [19, 20]. Wager and Athey [20] describes honest trees as those for which any particular training example  $(X_i, Y_i)$  is used to partition feature space, or to estimate the quantity of interest, but not both. This property can be achieved in (at least) two ways. The first method is splitting the observed sample into two sets, one set for learning the partitions of feature space  $\mathcal{X}$ , and one set to make vote on the plurality or average within each leaf node. We refer to them as the "partition" set  $\mathcal{D}^{\mathsf{P}}$  and "voting" set  $\mathcal{D}^{\mathsf{V}}$ , respectively. (Denil et al. [3] called them "structure" and "estimation", which we avoid because both sets are used to estimate different quantities).

For example, say we wish to estimate the conditional mean function  $\mu(x) = \mathbb{E}[Y \mid X = x]$ . Let  $\mathcal{L}$  be a partition of feature space  $\mathcal{X}$ , as described in Section 2.3. Letting m < n, such an  $\mathcal{L}$  can be learned via a decision tree with  $\mathcal{D}^{\mathsf{P}} = \{(X_1, Y_1), ..., (X_m, Y_m)\}$ . This leaves  $\mathcal{D}^{\mathsf{V}} = \{(X_{m+1}, Y_{m+1}), ..., (X_n, Y_n)\}$ . Letting L(x) be the part of  $\mathcal{L}$  to which x belongs, the conditional mean estimate can be  $\hat{\mu}(x) = \frac{\sum_{\mathcal{D}^{\mathsf{V}}} Y_i \cdot \mathbb{I}[X_i \in L(x)]}{\sum_{\mathcal{D}^{\mathsf{V}}} \mathbb{I}[X_i \in L(x)]} = \frac{\sum_{i=m+1}^n Y_i \cdot \mathbb{I}[X_i \in L(x)]}{\sum_{i=m+1}^n \mathbb{I}[X_i \in L(x)]}$ . This method can be applied at the forest level, that is, using the same partition points to learn every decision tree, or at the tree level, in which the data is randomly partitioned into partition and voting sets in each tree. Denil et al. [3] has shown that tree level splitting increases performance for low sample sizes, while for higher sample sizes the performance difference is indistinguishable. A second way of achieving honesty is by ignoring the responses of interest  $Y_i$  and only using  $X_i$  and any auxiliary variables  $W_i$  to place splits in each decision tree. Wager and Athey [20] proposes an example of this method in estimating heterogeneous treatment effects via random forest.

Balanced trees are described by Wager and Athey [20] and Athey et al. [19] as  $\alpha$ -regular, meaning that at each split in a tree, at least an  $\alpha$  fraction of the data are placed in each child node. Finally, trees for which each dimension has a nonzero chance of being used for the split are called random-split. This can be achieved, for example, by randomly choosing one candidate dimension uniformly at each split. For honest, balanced, random-split trees, consistency can be shown for a very general class of random forest estimates [19].

**3 Uncertainty Forests** Given observations  $\mathcal{D}_n = \{(X_1, Y_1), ..., (X_n, Y_n)\}$ , the goal is to estimate conditional entropy  $H(Y|X) = \mathbb{E}_{X'}[H(Y|X=X')]$  and use that to estimate mutual information. UF provides a consistent estimate of the conditional entropy, and an empirically non-biased estimate with sufficient sample size.

UF partitions the data into three sets: a partition set  $\mathcal{D}^{\mathsf{P}}$ , a voting set  $\mathcal{D}^{\mathsf{V}}$ , and evaluation set  $\mathcal{D}^{\mathsf{E}}$ .  $\mathcal{D}^{\mathsf{P}}$  and  $\mathcal{D}^{\mathsf{V}}$  are used to learn low-bias posteriors  $\hat{p}(y\mid x)$ , and consequently the function  $\hat{H}(Y\mid X=x)$ .  $\mathcal{D}^{\mathsf{E}}$  will be used to estimate its expectation  $\mathbb{E}_{X'}[H(Y\mid X=X')]$ . Note that for forest-level evaluation,  $\mathcal{D}^{\mathsf{E}}$  need not be labeled, as only the x values are necessary to evaluate the function  $H(Y\mid X=x)$ . If data is difficult to label, UF can effectively leverage unlabelled data in a semisupervised fashion.

#### 3.1 Algorithm

0. Set hyper-parameters including tree convergence criteria, number of trees B, and finite sample correction coefficient  $\kappa$ , and (possibly unlabelled) evaluation set  $\mathcal{D}^{\mathsf{E}}$ .

### 1. Build Trees

For each tree b from 1 to B (the maximum number of trees):

- (a) Randomly subsample  $s \leq n$  data points from  $\mathcal{D}_n \mathcal{D}^{\mathsf{E}}$ .
- (b) Randomly split the s data points into a partition set  $\mathcal{D}_b^{\mathsf{P}}$  and voting set  $\mathcal{D}_b^{\mathsf{V}}$ .
- (c) Using  $\mathcal{D}_b^{\mathsf{P}}$ , learn a decision tree partition  $\mathcal{L}_b$ .  $L_b(x)$  is the leaf node of  $\mathcal{L}_b$  that x "falls" into.
- (d) Using the  $\mathcal{D}_b^{\mathsf{V}}$ , letting  $N_b(x) = \sum_{i \in \mathcal{D}^{\mathsf{V}}} \mathbb{1}[X_i \in L_b(x)]$  be the leaf size of  $L_b(x)$ , estimate

the conditional probability of Y=y given X=x by  $\tilde{p}_b(y\mid x)=\frac{1}{N_b(x)}\sum_{i\in\mathcal{D}^{\mathsf{V}}}\mathbb{1}[Y_i=y \text{ and } X_i\in L_b(x)].$  This is the empirical frequency of y (given by the voting data) in the leaf node of x.

(e) When Y is categorical, all samples in a leaf estimator may belong to one class even though the probabilities for other classes are nonzero, which biases finite-sample estimates of the conditional distribution. UF remedies this by adapting a robust finite sampling technique described in Vogelstein et al. [21]. We first replace all zero probabilities with  $1/(\kappa N_b(x))$  (where  $\kappa>0$  is some suitably chosen constant), and renormalize the probabilities. The finite-sample corrected estimate is thus  $\hat{p}_b(y\mid x)=\frac{\hat{p}_b(y\mid x)}{\sum_{k=1}^K \hat{p}_b(k\mid x)}$ . If  $N_b(x)\stackrel{P}{\to}\infty$ , then  $\hat{p}_b(y\mid x)$  approaches  $\tilde{p}_b(y\mid x)$ .

# 2. Estimate conditional entropy function

- (a) Average all of the posterior estimates from each tree:  $\hat{p}(y \mid x) = \frac{1}{B} \sum_{b=1}^{B} \hat{p}_b(y \mid x)$ .
- (b) Set the conditional entropy function estimator  $\hat{H}(Y \mid X = x) = -\sum_{y \in [K]} \hat{p}(y \mid x) \log \hat{p}(y \mid x)$ .

# 3. Compute conditional entropy and mutual information

- (a) Evaluate at every point in  $\mathcal{D}^{\mathsf{E}}$  and average to yield  $\hat{H}(Y\mid X)=\frac{1}{|\mathcal{D}^{\mathsf{E}}|}\sum_{i\in\mathcal{D}^{\mathsf{E}}}\hat{H}(Y\mid X=X_i)$ .
- (b) Letting  $\hat{p}(y) = \frac{1}{n} \sum_{i \in \mathcal{D}_n} \mathbb{1}[Y_i = y]$ , estimate H(Y) with  $\hat{H}(Y) = -\sum_{y \in [K]} \hat{p}(y) \log \hat{p}(y)$ .
- (c) Let  $\hat{I}(X;Y) = \hat{H}(Y) \hat{H}(Y \mid X)$ .

UF uses basic CART for tree construction, which has a number of hyper-parameters. These include which objective function to optimize for each split, how many samples to use to generate each tree (s), minimum samples in a leaf before splitting (k), maximum tree depth, and number of trees (B). Random forest has been shown in practice to be very robust to hyperparameters [22]. For UF, we do not impose a maximum depth and use general rules-of-thumb for the other choices (minimize Gini impurity, k=1, B=300). We use  $|\mathcal{D}^{\mathsf{P}}|=0.4\cdot n$ ,  $|\mathcal{D}^{\mathsf{V}}|=|\mathcal{D}^{\mathsf{E}}|=0.3\cdot n$  in experiments  $(s=|\mathcal{D}^{\mathsf{P}}|+|\mathcal{D}^{\mathsf{V}}|)$ . The UF pseudocode is described in detail in the supplementary material. In experiments, we find that letting the evaluation set  $\mathcal{D}^{\mathsf{E}}$  be randomly sampled at the tree-level performs better than holding out a set that is common for all trees. We also use  $\kappa=3$ , after sweeping over choices in [0.1,10]. Letting T be the number of threads, the learning time complexity is  $O(Bdn\log^2(n))$ , the memory complexity is O(B(nd+K)), and the storage complexity is O(Bn). [23]

**3.2** Consistency of Uncertainty Forest Estimates Under mild conditions, UF provides a consistent estimate of conditional probability, conditional entropy, and mutual information. All proofs are collected in the supplementary material. Assume that  $\mathcal{Y}$  is discrete  $(\mathcal{Y}=[K])$ , and that  $\mathcal{X}=\mathbb{R}^d$ . Let  $\hat{H}_n(Y\mid X)$  be the conditional entropy estimate and  $\hat{I}_n(X;Y)$  be the mutual information estimate, now indexed by n for explicitness. Suppose that  $|\mathcal{D}^{\mathsf{E}}|=\gamma n$  for  $0<\gamma<1$ . Suppose that the subsample size  $s_n$  is chosen such that  $s_n\to\infty$  and  $\frac{s_n}{n}\to0$ . Consider the following assumptions.

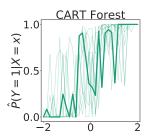
Assumption 1. Suppose that X is supported on  $\mathbb{R}^d$  and has a density which non-zero almost everywhere. This is equivalent to have X be uniformly distributed in  $[0,1]^d$  due to the monotone invariance of trees [3].

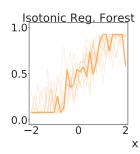
Assumption 2. Suppose that for each  $y \in [K]$  the conditional probability  $p(y \mid x)$  is Lipshitz continuous on  $\mathbb{R}^d$ .

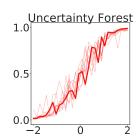
Assumption 3. Each tree b is constructed such that for all  $x \in \mathcal{X}, \epsilon > 0$ ,  $P[N_b(x) < \epsilon] \to 0$  as  $n \to \infty$ .

Theorem 1 (Consistency of the conditional entropy estimate). Given Assumptions 1-3, the conditional entropy estimate is consistent as  $n \to \infty$ , that is,  $\hat{H}_n(Y \mid X) \stackrel{P}{\to} H(Y \mid X)$ .

This result states that the estimate  $\hat{H}(Y \mid X)$  is arbitrarily close in probability to the true  $H(Y \mid X)$  for sufficiently large n. As a simple consequence, due to the consistency of the maximum-likelihood







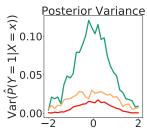


Figure 1: Comparison of estimated posterior distributions using random forest algorithms. Left plots show posterior distribution of Y=1 given x from CART, honest, and UF. Ten trials are plotted for each algorithm, the mean is highlighted. Right-most plot shows variance, over 100 trials, of posterior estimates vs x.  $\mu=1, n=6000$  for all plots.

estimate  $\hat{H}(Y)$  for H(Y) based on empirical frequencies, we have Theorem 2. Both results rely heavily on recent theoretical work by Athey et al. [19].

Theorem 2 (Consistency of the mutual information estimate). With the same assumptions as in Theorem 1,  $\hat{I}_n(X;Y) \stackrel{P}{\to} I(X;Y)$  as  $n \to \infty$ .

#### 4 Simulation Results

**4.1 Posterior Probability Experiments** Consider the following setting: let each  $Y_i$  be Bernoulli with 50% probability to be either +1 or -1; let each  $X_i$  be normally distributed with mean  $Y_i \times \mu$  and variance one, where  $\mu$  is a parameter controlling effect size. A CART random forest, a CART forest with isotonic regression calibration [14], and UF (using both honest sampling and finite sample correction) are trained on data drawn from the above distribution with  $\mu=1$ . We estimate posterior distributions and plot the posterior for class 1 in Figure 1. As x increases, the true probability that Y is one givne X=x increases. Thus, unsurprisingly, all random forest algorithms have  $\hat{P}(Y=1 \mid X=x)$  decrease to 0 as x becomes more negative, and increase to 1 as x becomes more positive. However, the posterior estimated from UF has significantly lower variance than both normal CART forests and isotonic regression forests (Figure 1, right).

**4.2 Conditional Entropy Experiments** These better posterior estimates from UF carry over to better estimates of conditional entropy. Figure 2A shows that UF estimates converge to the truth as sample size increases, whereas isotonic regression forest estimates and CART forest estimates are biased even for large sample sizes in the low-dimensional setting. All three algorithms behave as expected when the effect size  $(\mu)$  increases: they all approach zero conditional entropy, as seen in Figure 2B. CART forests and isotonic regression forest however, exhibit a bias for small effect size. For the high-dimensional experiment,  $X_i$ 's are multivariate Gaussians, where the mean of the first dimension is still  $Y_i \times \mu$  but each additional dimension has mean 0 and identity covariance:  $X_i \sim \mathcal{N}(\underbrace{(Y_i\mu,0,\ldots,0)^T}_{d},\underbrace{I_d}_{ldentity\ matrix})$ , where  $\mathcal{N}(\theta,\Sigma)$  refers to the Gaussian distribution with mean

 $\theta \in \mathbb{R}^d$  and covariance matrix  $\Sigma \in \mathbb{R}^{d \times d}$ . Because each added dimension is noise, the conditional entropy does not change. This allows us to compare behavior of our forest estimates to truth [24]. Figure 2C show that when d=20, the UF estimate still converges to truth as sample size increases. Interestingly, the bias of isotonic regression forests decreases in the high-dimensional setting, suggesting that in this setting the true posteriors might be (approximately) monotonic in the learned posteriors. When d=1, the bias of IRF becomes worse as sample size increases. Figure 2D demonstrates that CART forests remain highly biased, whereas the other approaches lack bias in these high-dimensional settings.

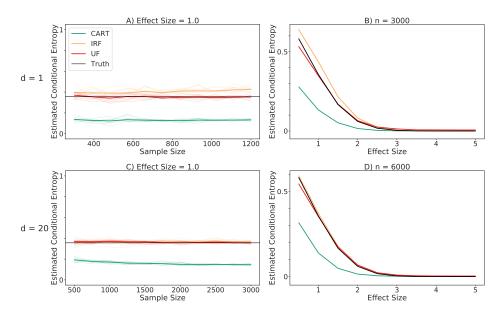


Figure 2: Behavior of random forest estimates for conditional entropy. Top plots are for d=1; bottom plots are for d=20. The left plot shows estimates vs. increasing sample size ( $\mu=1$ ). Twenty trials are plotted with high transparency to show variance. Right plot shows estimates vs. increasing  $\mu$  (n=3000 for d=1 and n=6000 for d=20). UF is consistent, and other approaches remain biased for very large sample sizes. UF is also closest to the truth for low effect sizes.

**4.3 Mutual Information Experiments** We compare UF to the KSG, mixed KSG, and isotonic regression estimators of mutual information [4, 9, 14]. Consider three simulation settings, all based on mixtures of Gaussians with various parameters. We compute normalized mutual information,  $I(X;Y)/\min\{H(X),H(Y)\}$ . For each setting, we consider both dimensionality d=1 and increasing dimension up to d=20. Only the first dimension depends on the class label, each additional dimension is an independent, standard Gaussian. Because only the first dimension contains any signal, mutual information does not change with increasing dimensionality [24]. In the case of two classes, Y is drawn Bernoulli with probability  $\pi$ . For the three class case, the classes are drawn multinomial from the vector  $(\pi, \frac{1-\pi}{2}, \frac{1-\pi}{2})$ . We explore the robustness of the estimators both with d=2 and changing class prior  $\pi$ , as well as increasing dimensionality when  $\pi=\frac{1}{2}$ .

**Spherical Gaussians** A mixture of two Gaussians, the same distribution as in Figure 2. The  $\mu$  parameter controls the effect size while  $y \in \{-1, +1\}$  controls the class label:  $X \mid Y = y \sim \mathcal{N}(y(\mu, 0)^\mathsf{T}, I)$ .

**Elliptical Gaussians** To quantify performance of MI estimators when the Bayes optimal decision boundary is not axis-aligned, consider elliptical Gaussians:  $X \mid Y = y \sim \mathcal{N}\left(y(\mu,0)^\mathsf{T}, \Sigma_y\right)$ , where  $\Sigma_{-1} = \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix}$ , and  $\Sigma_{+1} = I$ .

**Three Classes** To quantify performance for greater than two classes, let  $y \in \{0, 1, 2\}$ , and  $X \mid Y = u \sim \mathcal{N}(\boldsymbol{\mu}_n, I)$  where  $\boldsymbol{\mu}_0 = (0, u)^\mathsf{T}$ ,  $\boldsymbol{\mu}_1 = (u, 0)^\mathsf{T}$ ,  $\boldsymbol{\mu}_2 = (-u, 0)^\mathsf{T}$ .

 $y \sim \mathcal{N}(\pmb{\mu}_y, I)$  where  $\pmb{\mu}_0 = (0, \mu)^\mathsf{T}, \;\; \pmb{\mu}_1 = (\mu, 0)^\mathsf{T}, I), \;\; \pmb{\mu}_2 = (-\mu, 0)^\mathsf{T}.$  Figure 3 shows the performance of each estimator. When d=2, UF, KSG, mixed KSG, and IRF all do reasonably well, though IRF is heavily biased for various level of class balance. As dimension increases, the KSG and Mixed KSG estimators suffer a significant performance degradation. In the three class case, the KSG suffers a worse bias in high-dimensional settings. On the other hand, the UF estimator maintains performance as dimensionality gets high, but incurs more bias amid severe imbalance. For these Gaussian settings, IRF actually improves its performance from low-to-high dimensional data. In comparison to other methods, UF shows strong performance in high-dimensional and moderately imbalanced settings.

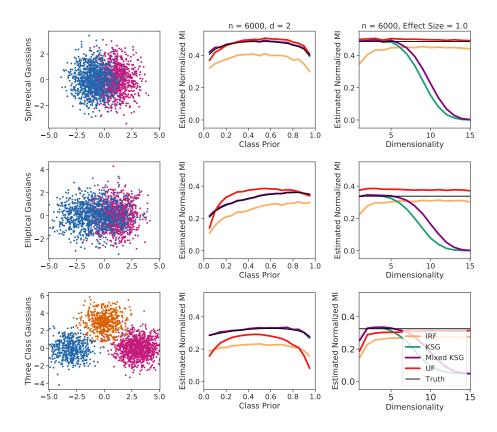


Figure 3: Mutual information estimates in three different ("near") Gaussian settings. Left: an example sample for each setting. Center. Normalized mutual information for d=2 and n=6000. Right. Normalized mutual information estimates at n=6000 for various dimensions. The mean of twenty trials is plotted. Both KSG and Mixed KSG break down in the high-dimensional setting, past d=8. Similar to 2, isotonic regression is biased in the low-dimensional setting, but improves for these distributions for high dimensions. For the two-class settings, UF successfully recovers the mutual information, while is suffers slight bias in the three-class setting. See main text for equations.

Mutual Information in Drosophila Neural Network (Connectome) An immediate application of our random forest estimate of conditional entropy is measuring information contained in neuron types for the larval Drosophila mushroom body (MB) connectome [25]. This dataset, obtained via serial section transmission electron microscopy, provides a real and important opportunity for investigating synapse-level structural connectome modeling [26]. This connectome consists of 213 different neurons (n = 213) in four distinct cell types: Kenyon Cells, Input Neurons, Output Neurons, and Projection Neurons. The connectome adjacency matrix is visualized in Figure 4 (left). Each neuron comes with a mixture of categorical and continuous features: "claw" refers to the integer number of dendritic claws for Kenyon cells, "dist" refers to real distance from the neuron to the neuropil, "age" refers to neuron (normalized) age as a real number between -1 and 1, and "cluster" refers to the community detected by the latent structure model as in Priebe et al. [27]. We compute mutual information with Y as the neuron type and X as various subsets of the features. Because neuron type has been a subjective categorical assignment based on gross morphological features, we expect mutual information to be high for the entire feature vector. Running a permutation test with 1000 replicates, the test statistic  $\hat{I}_n(Y,X) = 0.913$  was found to be statistically greater than zero with a p-value of 0.001. Regarding the individual features, scientific prior knowledge posits a few relationships that are confirmed by the mutual information estimates. Letting  $X_{in}$  be the subset of features under consideration, we compute  $\hat{I}(Y, X_{in})$ , the mutual information between Y and the "in" features, and  $\hat{I}(Y, X_{\text{out}} \mid X_{\text{in}})$ , the additional information given by the "out" features (Figure 4 (right)). We confirm that the estimates  $\ddot{I}(Y, X_{\text{in}}) + \ddot{I}(Y, X_{\text{out}} \mid X_{\text{in}}) =$ 

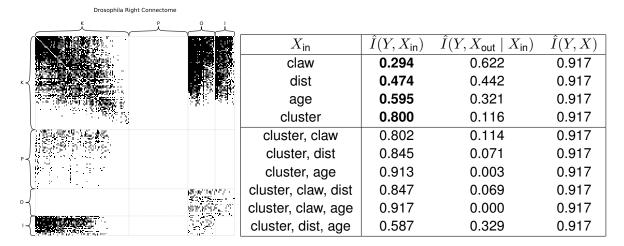


Figure 4: Left Drosophila larva right hemisphere connectome. Groups of neurons are labelled K for Kenyon Cells, I for Input Neurons, O for Output Neurons, and P for Projection Neurons (PN). Black cells represent the presence of an edge between the two corresponding nodes. Right Adjacency spectral embedding applied to the MB connectome shows clear cluster groups for each neuron type. This suggests a strong dependency between neuron type and neural features.

 $\hat{I}(X,Y)$  adhering to the chain rule of mutual information. Because only Kenyan cells have dendritic claws, this feature is unable to discriminate between the other three classes, explaining why it has the lowest mutual information estimate among the features. Age is typically computed using distance from the neuropil, which explains why both of these features yielded similar information regarding Y. Figure 4 of Priebe et al. [27] presents compelling evidence of latent structure model clusters being closely related to cell type, which is corroborated by its highest mutual information estimate among neuron features. Moreover, adding one or two additional features increases our estimated MI with respect to neuron class.

6 Conclusion We present Uncertainty Forest (UF), a nonparametric method of consistently estimating conditional entropy through randomized decision trees. Empirically, UF performs well in low- and high-dimensional settings. Furthermore, when extending our estimator to estimate mutual information, UF performs better than the mixed KSG and KSG estimators in a variety of settings. UF has strong theoretical justification in comparison to calibration methods such as isotonic regression forests. In machine learning and statistics, this tool could be adapted for feature selection, while in biology and neuroscience, scientists can find quantify uncertainty between various properties and labels.

The main limitation of this work is that it is only able to estimate uncertainty quantities for categorical Y; that said, the UF algorithm can be modified for continuous Y as well. Computing the posterior distribution  $\hat{P}(Y|X=x)$  when Y is continuous can be accomplished with a kernel density estimate instead of simply binning the probabilities. A theoretical and empirical analysis of this extension is of interest. When Y is multivariate, a heuristic approach such as subsampling Y dimensions or using multivariate random forests can be explored.

On the theoretical side, important next steps include rigorous proofs for convergence rates. Studying the behavior of UF estimates in more complicated nonlinear, high-dimensional settings should be explored as well. Practical applications such as dependence testing and k-sample testing for high-dimensional, nonlinear data will be natural applications for these information theoretic estimates.

Data and Code Availability Statement The implementation of UF, the simulated and real data experiments, and their visualization can be reproduced completely with the instructions in code available on GitHub. The experiments were run in parallel on a 1TB RAM machine with 96-cores.

**Acknowledgements** The authors are grateful for the support by the XDATA program of the Defense Advanced Research Projects Agency (DARPA) administered through Air Force Research Laboratory contract FA8750-12-2-0303, and DARPA's Lifelong Learning Machines program through contract FA8650-18-2-7834, and the National Science Foundation award DMS-1921310. Research was partially supported by funding from Microsoft Research.

#### References

- [1] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 531–540, StockholmsmÃd'ssan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/belghazi18a.
- [2] N. Meinshausen. Quantile regression forests. *Journal of Machine Learning Research*, 7:983–999, 2006.
- [3] M. Denil, D. Matheson, and N. De Freitas. Narrowing the gap: Random forests in theory and in practice. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 665–673, Jun 2014
- [4] Weihao Gao, Sreeram Kannan, Sewoong Oh, and Pramod Viswanath. Estimating mutual information for discrete-continuous mixtures. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5986–5997. Curran Associates, Inc., 2017. URL http://papers.nips.cc/paper/7180-estimating-mutual-information-for-discrete-continuous-mixtures.pdf.
- [5] J Beirlant, E J Dudewicz, L Györfi, and E C Van Der Meulen. Nonparametric entropy estimation: An overview. *International Journal of Mathematical and Statistical Sciences*, 6(1): 17–39, 1997. URL http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.87.5281&rep=rep1&type=pdf.
- [6] Nikolai Leonenko, Luc Pronzato, and Vippal Savani. Estimation of entropies and divergences via nearest neighbors, 2008.
- [7] Thomas B. Berrett, Richard J. Samworth, and Ming Yuan. Efficient multivariate entropy estimation via k-nearest neighbour distances. *Ann. Statist.*, 47(1):288–318, 02 2019. doi: 10.1214/18-AOS1688. URL https://doi.org/10.1214/18-AOS1688.
- [8] K. Sricharan, D. Wei, and A. O. Hero. Ensemble estimators for multivariate entropy estimation. *IEEE Trans Inf Theory*, 59(7):4374–4388, Jul 2013.
- [9] Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. Phys. Rev. E, 69:066138, Jun 2004. doi: 10.1103/PhysRevE.69.066138. URL https://link.aps.org/doi/10.1103/PhysRevE.69.066138.
- [10] Andrew M. Fraser and Harry L. Swinney. Independent coordinates for strange attractors from mutual information. *Phys. Rev. A*, 33:1134–1140, Feb 1986. doi: 10.1103/PhysRevA.33.1134. URL https://link.aps.org/doi/10.1103/PhysRevA.33.1134.
- [11] Kirthevasan Kandasamy, Akshay Krishnamurthy, Barnabas Poczos, Larry Wasserman, and james m robins. Nonparametric von mises estimators for entropies, divergences and mutual informations. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, Advances in Neural Information Processing Systems 28, pages 397–405. Curran Associates, Inc., 2015. URL <a href="http://papers.nips.cc/paper/5911-nonparametric-von-mises-estimators-for-entropies-divergences-and-mutual-informations.">http://papers.nips.cc/paper/5911-nonparametric-von-mises-estimators-for-entropies-divergences-and-mutual-informations.pdf.</a>
- [12] Shuyang Gao, Greg Ver Steeg, and Aram Galstyan. Efficient estimation of mutual information for strongly dependent variables. *CoRR*, abs/1411.2003, 2014. URL <a href="http://arxiv.org/abs/1411.2003">http://arxiv.org/abs/1411.2003</a>.

- [13] Marylou Gabrié, Andre Manoel, Clément Luneau, Jean Barbier, Nicolas Macris, Florent Krzakala, and Lenka Zdeborová. Entropy and mutual information in models of deep neural networks. ArXiv, abs/1805.09785, 2018.
- [14] Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD âĂŹ02, page 694âĂŞ699, New York, NY, USA, 2002. Association for Computing Machinery. ISBN 158113567X. doi: 10.1145/775047.775151. URL https://doi.org/10.1145/775047.775151.
- [15] L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001. ISSN 1573-0565.
- [16] Manuel Fernández-Delgado, Eva Cernadas, Senén Barro, and Dinani Amorim. Do we need hundreds of classifiers to solve real world classification problems? *Journal of Machine Learning Research*, 15:3133–3181, 2014.
- [17] Bingguo Li, Xiaojun Chen, Mark Junjie Li, Joshua Zhexue Huang, and Shengzhong Feng. Scalable random forests for massive data. In *Proceedings of the 16th Pacific-Asia Conference on Advances in Knowledge Discovery and Data Mining Volume Part I*, PAKDD'12, pages 135–146, Berlin, Heidelberg, 2012. Springer-Verlag. ISBN 978-3-642-30216-9. doi: 10.1007/978-3-642-30217-6\_12. URL http://dx.doi.org/10.1007/978-3-642-30217-6\_12.
- [18] Leo Breiman. Classification and Regression Trees. CRC Press, hardcover edition, 2017. ISBN 1138469521,9781138469525. URL http://gen.lib.rus.ec/book/index.php?md5= 1D398E5BB3FA0334744196C49276052D.
- [19] Susan Athey, Julie Tibshirani, and Stefan Wager. Generalized random forests. *Ann. Stat.*, 47(2): 1148–1178, April 2019. URL https://projecteuclid.org/euclid.aos/1547197251.
- [20] Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018. doi: 10.1080/01621459.2017.1319839. URL https://doi.org/10.1080/01621459.2017.1319839.
- [21] J. T. Vogelstein, W. Gray Roncal, R. J. Vogelstein, and C. E. Priebe. Graph classification using signal-subgraphs: Applications in statistical connectomics. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 35(7):1539–1551, July 2013. ISSN 0162-8828.
- [22] Philipp Probst, Anne-Laure Boulesteix, and Bernd Bischl. Tunability: Importance of hyperparameters of machine learning algorithms. *Journal of Machine Learning Research*, 20(53):1–32, 2019. URL http://jmlr.org/papers/v20/18-444.html.
- [23] Tyler M. Tomita, James Browne, Cencheng Shen, Jaewon Chung, Jesse L. Patsolic, Benjamin Falk, Jason Yim, Carey E. Priebe, Randal Burns, Mauro Maggioni, and Joshua T. Vogelstein. Sparse projection oblique randomer forests, 2015.
- [24] Aaditya Ramdas, Sashank J. Reddi, Barnabás Póczos, Aarti Singh, and Larry Wasserman. On the decreasing power of kernel and distance based nonparametric hypothesis tests in high dimensions. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAIâĂŹ15, page 3571âĂŞ3577. AAAI Press, 2015. ISBN 0262511290.
- [25] K. Eichler, F. Li, A. Litwin-Kumar, Y. Park, I. Andrade, C. M. Schneider-Mizell, T. Saumweber, A. Huser, C. Eschbach, B. Gerber, R. D. Fetter, J. W. Truman, C. E. Priebe, L. F. Abbott, A. S. Thum, M. Zlatic, and A. Cardona. The complete connectome of a learning and memory centre in an insect brain. *Nature*, 548(7666):175–182, 08 2017.
- [26] Joshua T Vogelstein, Eric W Bridgeford, Benjamin D Pedigo, Jaewon Chung, Keith Levin, Brett Mensh, and Carey E Priebe. Connectal coding: discovering the structures linking cognitive phenotypes to individual histories. *Curr. Opin. Neurobiol.*, 55:199–212, May 2019. URL <a href="http://dx.doi.org/10.1016/j.conb.2019.04.005">http://dx.doi.org/10.1016/j.conb.2019.04.005</a>.
- [27] Carey E. Priebe, Youngser Park, Minh Tang, Avanti Athreya, Vince Lyzinski, Joshua T. Vogelstein, Yichen Qin, Ben Cocanougher, Katharina Eichler, Marta Zlatic, and Albert Cardona. Semiparametric spectral modeling of the drosophila connectome, 2017.
- [28] P. Langley. Crafting papers on machine learning. In Pat Langley, editor, Proceedings of the 17th

- International Conference on Machine Learning (ICML 2000), pages 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- [29] Leo Breiman. Classification and Regression Trees. Routledge, 1984.
- [30] S. Athey, J. Tibshirani, and S. Wager. Generalized random forests. *Annals of Statistics*, 47(2): 1148–1178, 2019.
- [31] J. T. Vogelstein, Q. Wang, E. Bridgeford, C. E. Priebe, M. Maggioni, and C. Shen. Discovering and deciphering relationships across disparate data modalities. *eLife*, 8:e41690, 2019. doi: 10.7554/ elife.41690.
- [32] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- [33] J. H. Zhang, T. D. Chung, and K. R. Oldenburg. A simple statistical parameter for use in evaluation and validation of high throughput screening assays. *J Biomol Screen*, 4(2):67–73, 1999.
- [34] J. W. Prescott. Auantitative imaging biomarkers: the application of advanced image processing and analysis to clinical and preclinical decision making. *J Digit Imaging*, 26(1):97–108, February 2013.
- [35] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, NY, USA, 2000. ISBN 0-521-77362-8.
- [36] B. W. Silverman. Weak and strong uniform consistency of the kernel estimate of a density and its derivatives. *Annals of Statistics*, 6(1):177–184, Jan 1978.
- [37] Kenneth Wallis. A note on the calculation of entropy from histograms. MPRA Paper 52856, University Library of Munich, Germany, October 2006. URL https://ideas.repec.org/p/pra/mprapa/52856. html.
- [38] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

Appendix A. Pseudocode. The FITRANDOMFOREST, GETINBAGSAMPLES, and APPLYTREE operations, as well as the NUMCLASSES and NUMLEAVES fields are all standard functions in the scikit-learn decision tree and bagging classifier modules [38].

## Algorithm 1 Uncertainty Forest (Forest-level)

```
Input: Training set \mathcal{D}_n and evaluation set \mathcal{D}^{\mathsf{E}}, \theta = \{\text{number of trees } B, \text{minimum leaf size } k, \text{subsample } k\}
   size s }, finite correction constant \kappa.
Output: Conditional entropy estimate \hat{H}(Y \mid X).
   function UncertaintyForest(\mathcal{D}_n, \mathcal{D}^{\mathsf{E}}, \theta, \kappa)
        \mathcal{T} = \mathcal{D}_n - \mathcal{D}^{\mathsf{E}}
        model = FITRANDOMFOREST(T, \theta).
        posteriors = [].
        for tree b in model do
            posterior = ESTIMATEPOSTERIOR(\mathcal{T}, model, \kappa, b)
            posteriors.append(posterior)
        \hat{H}(Y \mid X) = \text{EVALUATEPOSTERIORS}(\mathcal{D}^{\mathsf{E}}, \text{ model, posteriors}).
        return \hat{H}(Y \mid X).
   end function
```

```
Algorithm 2 Uncertainty Forest (Tree-level)
Input: Training set \mathcal{D}_n and evaluation set \mathcal{D}^{\mathsf{E}}, \theta = \{\text{number of trees } B, \text{ minimum leaf size } k, \text{ subsample } k\}
   size s }, finite correction constant \kappa.
Output: Conditional entropy estimate \hat{H}(Y \mid X).
   function UNCERTAINTYFOREST(\mathcal{D}_n, \mathcal{D}^{\mathsf{E}}, \theta, \kappa)
       \mathcal{T} = \mathcal{D}_n - \mathcal{D}^{\mathsf{E}}
       model = FITRANDOMFOREST(\mathcal{T}, \theta).
       conditional_entropies = []
       for tree b in model do
            posterior = ESTIMATEPOSTERIOR(\mathcal{T}, model, \kappa, b)
            conditional_entropies.append(EVALUATEPOSTERIORS(\mathcal{D}^{E}, [model[b]], [posterior])).
       end for
       \ddot{H}(Y \mid X) = \text{conditional entropies.MEAN()}.
       return \hat{H}(Y \mid X).
   end function
```

Note that  $\mathcal{D}^{\mathsf{E}}$  can contain unlabelled data, as only the  $x_i$  values are used.

# **Algorithm 3** Posterior Estimation

```
Input: training set \mathcal{T} (full dataset without evaluation), fitted random forest model, finite correction con-
  stant \kappa, tree index b
Output: Posterior probability estimate tree b.
  function EstimatePosterior(\mathcal{T}, model, \kappa, b)
      K = \mathsf{model}.\mathsf{NumClasses}.
      \mathcal{D}_{P} = GETINBAGSAMPLES(model, b)
      \mathcal{D}_V = \mathcal{T} - \mathcal{D}_P.
      L = model[b].NumLeaves.
      vote_counts = [0]^{L \times K}.
      for observation (x, y) in \mathcal{D}_V do
          l = model[b].ApplyTree(x)
          vote\_counts[l, y] = vote\_counts[l, y] + 1.
      end for
      leaf sizes = RowSum(vote counts).
      posterior = [0]^{L \times K}.
      for leaf index l in [L] do
          for class y in [K] do
              posterior[l, y] = vote counts[l, y] / leaf sizes[l].
          end for
      end for
      posterior = FINITESAMPLECORRECT(posterior, \kappa, leaf_sizes)
      return posterior
  end function
```

# Algorithm 4 Posterior Evaluation

```
Input: evaluation set \mathcal{D}^{E}, fitted random forest model, posteriors for each tree
Output: Conditional entropy estimate \hat{H}(Y \mid X).
   function EVALUATEPOSTERIORS(\mathcal{D}^{E}, model, posteriors)
        B = \mathsf{model}.\mathsf{NumTrees}.
        for observation x_i in \mathcal{D}^{\mathsf{E}} do
            for tree b in model do
                 l = \mathsf{model}[b].\mathsf{APPLYTREE}(x).
                 for class y in [K] do
                      posterior = posteriors[b].
                      p_b(y \mid x_i) = posterior[l, y].
                 end for
            end for
            p_i(y \mid x) = \frac{1}{B} \sum_{b=1}^{B} p_b(y \mid x_i).
            \hat{H}_i(Y \mid X) = -\sum_{y \in [K]} p_i(y \mid x) \log p_i(y \mid x).
        end for
        \hat{H}(Y\mid X) = \frac{1}{|\mathcal{D}^{\mathsf{E}}|} \sum_{i\in\mathcal{D}^{\mathsf{E}}} \hat{H}_i(Y\mid X).
        return \hat{H}(Y \mid X).
   end function
```

# Algorithm 5 Finite Sample Correction

```
Input: posterior, finite correction constant \kappa, leaf_sizes
Output: Corrected posterior.
  function FINITESAMPLECORRECT(posterior, \kappa, leaf_sizes)
      (L, K) = posterior.SHAPE.
      for node index l in [L] do
         for class y in [K] do
             if posterior[l, y] == 0 then
                 posterior[l, y] = \frac{1}{\kappa \cdot \text{leaf\_sizes}[l]}
             end if
         end for
      end for
      normalizing_factor = RowSum(posterior).
      for node index l in [L] do
         for class y in [K] do
             posterior[l, y] = posterior[l, y] / normalizing_factor[l]
         end for
      end for
      return posterior.
  end function
```

## Appendix B. Proofs.

In this section, we present consistency results regarding estimation of conditional entropy and mutual information via Uncertainty Forest. The argument follows as a nearly direct consequence of Athey et al. [19], in which random forests that are grown according to some specifications, and solve locally weighted estimating equations are consistent and asymptotically Gaussian. We review the assumptions.

Assumption 1. Suppose that X is supported on  $\mathbb{R}^d$  and has a density which non-zero almost everywhere. This is equivalent to have X be uniformly distributed in  $[0,1]^d$  due to the monotone invariance of trees [3].

Assumption 2. Suppose that for each  $y \in [K]$  the conditional probability  $p(y \mid x)$  is Lipshitz continuous on  $\mathbb{R}^d$ .

Assumption 3. Each tree b is constructed such that for all  $x \in \mathcal{X}, \epsilon > 0$ ,  $P[N_b(x) < \epsilon] \to 0$  as  $n \to \infty$ .

Next, we summarize the notation and main result of Athey et al. [19]. Let  $\theta(x)$  be a parameter that is implicitly defined as the solution to some equation  $M_{\theta}(x) = 0$ . For example, in estimating the conditional mean  $\theta(x) = \mathbb{E}[Y \mid X = x]$ , we can use

$$M_{\theta}(x) = \mathbb{E}[Y \mid X = x] - \theta(x) = 0.$$

To simplify notation, we suppress the dependency of x in  $\theta$  and simply write  $\theta = \theta(x)$  wherever this dependency can be deduced from the context. To develop a sample estimate for  $\theta$ , we start with a sample score function  $\psi_{\theta}$  such that

$$\mathbb{E}[\psi_{\theta}(Y) \mid X = x] = M_{\theta}(x).$$

Now, let us be given a dataset  $\{(X_1, Y_1), ..., (X_n, Y_n)\}$ . The sample estimate  $\hat{\theta}$  solves a locally weighted estimating equation

$$\sum_{i=1}^{n} \alpha_i(x) \psi_{\theta}(Y_i) = 0.$$

In the case of conditional mean estimation, we will have  $\psi_{\theta}(Y) = Y - \theta$ . The  $\alpha_i(x)$ 's weigh highly observations with  $X_i$  close to x, as to approximate the conditional expectation of  $\psi_{\theta}(Y)$ , or  $M_{\theta}(x)$ . In a random forest, these weights amount to being the empirical probability that the test point x shares a leaf with each training point  $X_i$ . Under mild assumptions, Athey et al. [19] have shown such a method to be consistent for the estimation of  $\theta(x)$ . For such a result to be used for our purpose of estimating conditional entropy, we define  $M_{\theta}(x)$  as

$$M_{\theta}(x) = p(y \mid x) - \theta$$

Using this definition, we can derive consistency of  $\hat{H}(Y \mid X = x)$ . Given this function of x, we have the conditional entropy written as

$$H(Y \mid X) = \mathbb{E}_{X'}[H(Y \mid X = X')]$$

We first confirm that the finite-sample corrected posterior  $\hat{p}_b(y \mid x)$  approaches the uncorrected  $\tilde{p}_b(y \mid x)$  for large n.

Lemma 1. As  $n \to \infty$ 

$$|\hat{p}(y\mid x) - \tilde{p}(y\mid x)| \stackrel{P}{\rightarrow} 0$$

*Proof.* We index the estimates with n to make explicit the dependence on sample size. For tree b, in the cases where  $0 < \tilde{p}_{b,n}(y \mid x) < 1$  for all y, we have  $|\hat{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)| = 0$ . Otherwise, for  $\mathcal{N} = \{y : \tilde{p}_{b,n}(y \mid x) = 0\}$  and  $y \in \mathcal{N}$ , then

$$|\hat{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)| = \frac{1}{\kappa N_{b,n}(x)}.$$

By assumption 3, we have that  $N_{b,n}(x) \stackrel{n}{\to} \infty$  for all x. Thus,

$$\frac{1}{\kappa N_{b,n}(x)} \stackrel{P}{\to} 0$$

Similarly, the constant  $c=\frac{|\mathcal{N}|}{\kappa N_{b,n}(x)}\leq \frac{K}{\kappa N_{b,n}(x)}\to 0$  in probability. Thus, for  $y\not\in\mathcal{N}$ ,

$$|\hat{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)| = |(1 - c)\tilde{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)|$$
  
=  $|c \cdot \tilde{p}_{b,n}(y \mid x)| \le c \to 0$ 

as  $n \to \infty$ . The estimate  $\hat{p}_n(y \mid x) = \frac{1}{B} \sum_{b=1}^B \hat{p}_{b,n}(y \mid x)$  is a finite average of the  $\hat{p}_{b,n}(y \mid x)$ . Given  $\epsilon > 0$ ,

$$P[|\hat{p}_{n}(y \mid x) - \tilde{p}_{n}(y \mid x)| > \epsilon] = P[|\frac{1}{B} \sum_{b=1}^{B} \hat{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)| > \epsilon]$$

$$\leq P[\sum_{b=1}^{B} |\hat{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)| > B\epsilon]$$

$$\leq P\left[\bigcup_{b=1}^{B} |\hat{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)| > \epsilon\right]$$

$$\leq \sum_{b=1}^{B} P[|\hat{p}_{b,n}(y \mid x) - \tilde{p}_{b,n}(y \mid x)| > \epsilon]$$

$$= B \cdot P[|\hat{p}_{1,n}(y \mid x) - \tilde{p}_{1,n}(y \mid x)| > \epsilon]$$

$$= B \cdot (P[|\hat{p}_{1,n}(y \mid x) - \tilde{p}_{1,n}(y \mid x)| > \epsilon \cap \tilde{p}_{1,n}(y \mid x) \neq 0]$$

$$+ P[|\hat{p}_{1,n}(y \mid x) - \tilde{p}_{1,n}(y \mid x)| > \epsilon \cap \tilde{p}_{1,n}(y \mid x) \neq 0])$$

$$\stackrel{n}{\to} 0$$

Second, we require that honestly estimated posteriors  $\tilde{p}_n(y \mid x)$  (without the finite sample correction) converge to the true posteriors. This can be done by application of Theorem 3 from Athey et al. [19].

Lemma 2. For all  $y \in [K]$  and  $x \in \mathbb{R}^d$ ,  $\tilde{p}_n(y \mid x) \stackrel{P}{\to} p(y \mid x)$  as  $n \to \infty$ .

*Proof.* For a fixed (x, y), we can define  $\theta$  as the solution to

$$M_{\theta}(x) = p(y \mid x) - \theta = 0.$$

The sample equivalent is thus

$$\psi_{\theta}(Y) = \mathbb{1}[Y = y] - \theta,$$

and we have that

$$\mathbb{E}[\psi_{\theta}(Y) \mid X = x] = M_{\theta}(x).$$

Because our forest is learned according to Specification 1 in Athey et al. [19], we must confirm Assumptions 1 - 6 from Section 3 of that paper to apply their result. Only Assumption 1 is a proper assumption on the distribution, whereas the remaining are confirmed true based on our choice of  $M_{\theta}$  and  $\psi_{\theta}$ .

- 1. For fixed  $\theta$ ,  $M_{\theta}(x) = p(y \mid x) \theta$  is Lipschitz in x. We took this as a standard assumption on the joint distribution of (X, Y).
- 2. For fixed x,  $M_{\theta}(x)$  is twice continuously differentiable in  $\theta$ ,  $\frac{\partial^2 M_{\theta}}{\partial \theta^2} = 0$  and  $\frac{\partial}{\partial \theta} M_{\theta} = -1$ , which is invertible.
- 3. The function  $\psi_{\theta}$  satisfies

$$\sup_{x \in \mathcal{X}} (\mathsf{Var}[\psi_{\theta}(Y) - \psi_{\theta}(Y)]) = \sup_{x \in \mathcal{X}} (\mathsf{Var}[\mathbb{1}[Y = y] - \theta - (\mathbb{1}[Y = y] - \theta')]) = \sup_{x \in \mathcal{X}} (\mathsf{Var}[\theta' - \theta']) = 0,$$

 $\text{ and hence } \sup_{x \in \mathcal{X}} (\mathsf{Var}[\psi_{\theta}(Y) - \psi_{\theta'}(Y)]) \leq L ||\theta - \theta'|| \text{ for all } \theta, \theta' \text{ and some } L \geq 0.$ 

- 4. The function  $\psi_{\theta}$  is itself Lipschitz in  $\theta$ .
- 5. The solution of

$$\sum_{i=1} \alpha_i(x)\psi_{\theta}(Y_i) = 0$$

in  $\theta$  exists, and is equal to

$$\hat{\theta}(x) = \sum_{i=1}^{n} \alpha_i(x) \mathbb{1}[Y_i = y]$$

6. The function  $\psi_{\theta}$  is the negative subgradient of a convex function, and  $M_{\theta}$  is the negative subgradient of a strongly convex function (with respect to  $\theta$ ). Choose

$$\Psi(\theta) = \frac{1}{2} (\mathbb{1}[Y = y] - \theta)^2,$$
  
$$\mathbf{M}(\theta) = \frac{1}{2} (p(y \mid x) - \theta)^2$$

as these functions.

Granted the above assumptions, by Theorem 3 of Athey et al. [19] we have for every x, y,

$$\hat{p}_n(y \mid x) \stackrel{P}{\to} p(y \mid x).$$

By Lemma 2 and continuity, we have the consistency of the honest forest estimate of conditional entropy. Let  $\tilde{H}(Y\mid X=x) = -\sum_{y\in [K]} \tilde{p}_n(y\mid x)\log \tilde{p}_n(y\mid x)$ .

Corollary 3. For each  $x \in \mathbb{R}^d$ ,

$$\tilde{H}(Y \mid X = x) \stackrel{P}{\to} H(Y \mid X = x).$$

**Proof.** The function

$$h(p) = \begin{cases} 0 & \text{if } p = 0 \\ -p \log p & \text{otherwise} \end{cases}$$

is continuous on [0,1]. Similarly, the finite sum  $\sum_{k=1}^K h(p_k)$  is continuous on  $\{(p_1,...,p_K): 0 \le p_k \le 1, \sum_{k=1}^K p_k = 1\}$ . Thus, by Lemma 2 and the continuous mapping theorem, we have the desired result.

Lemma 4. The Uncertainty Forest function estimate  $\hat{H}(Y \mid X = x)$  converges in probability to  $H(Y \mid X = x)$  for each  $x \in \mathbb{R}^d$ .

*Proof.* By continuity of  $\sum_{k=1}^K h(p_k)$  on  $\{(p_1,...,p_K): 0 \le p_k \le 1, \sum_{k=1}^K p_k = 1\}$ , and Lemma 1 we have that for any x,

(B.1) 
$$\left| \hat{H}(Y \mid X = x) - \tilde{H}(Y \mid X = x) \right| \stackrel{P}{\rightarrow} 0$$

Then, given  $\epsilon > 0$ ,

$$\begin{split} P[|\hat{H}(Y\mid X=x) - H(Y\mid X=x)| > \epsilon] &\leq P[|\hat{H}(Y\mid X=x) - \tilde{H}(Y\mid X=x)| \\ &+ |\tilde{H}(Y\mid X=x) - H(Y\mid X=x)| > \epsilon] \\ &\leq P[|\hat{H}(Y\mid X=x) - \tilde{H}(Y\mid X=x)| > \frac{\epsilon}{2}] \\ &+ P[|\tilde{H}(Y\mid X=x) - H(Y\mid X=x)| > \frac{\epsilon}{2}] \\ &\stackrel{n}{\to} 0, \end{split}$$

by Lemma 3 and B.1.

For simplicity of notation, let n refer only to the total sample size of the partition  $\mathcal{D}^{\mathsf{P}}$  and vote  $\mathcal{D}^{\mathsf{V}}$  sets. Let  $\nu$  be the size of the evaluation set  $\mathcal{D}^{\mathsf{E}}$ . If a subset of the labelled data is being used, then this may be some fraction  $0 < \gamma < 1$  of the size of the full training data. In this case, we consider any growing evaluation set size  $\nu$ .

Theorem 1 (Consistency of the conditional entropy estimate). Suppose the conditions of the preceding lemmas. Then the conditional entropy estimate is consistent as  $n, \nu \to \infty$ , that is,

$$\hat{H}_{n,\nu}(Y \mid X) \stackrel{P}{\to} H(Y \mid X).$$

*Proof of Theorem 1.* By Lemma 4 we have the pointwise consistency of  $\hat{H}_n(Y \mid X = x)$  (now indexed by n explicitly) for  $H(Y \mid X = x)$ . We compute the limits with respect to n and  $\nu$  in either order to achieve the result. Let  $\hat{H}_{n,\nu}(Y \mid X) = \frac{1}{\nu} \sum_{i \in \mathcal{D}^{\mathbb{E}}} \hat{H}_n(Y \mid X = X_i)$ .

$$\lim_{\nu \to \infty} \lim_{n \to \infty} \hat{H}_{n,\nu}(Y \mid X) = \lim_{\nu \to \infty} \frac{1}{\nu} \sum_{i \in \mathcal{D}^{\mathsf{E}}} \lim_{n \to \infty} \hat{H}_n(Y \mid X = X_i)$$

$$= \lim_{\nu \to \infty} \frac{1}{\nu} \sum_{i \in \mathcal{D}^{\mathsf{E}}} H(Y \mid X = X_i)$$

$$= H(Y \mid X)$$

by the Weak Law of Large Numbers, where the limits are taken in probability. On the other hand,

$$\lim_{n \to \infty} \lim_{\nu \to \infty} \hat{H}_{n,\nu}(Y \mid X) = \lim_{n \to \infty} \lim_{\nu \to \infty} \frac{1}{\nu} \sum_{i \in \mathcal{D}^{\mathsf{E}}} \hat{H}_n(Y \mid X = X_i)$$

$$= \lim_{n \to \infty} \mathbb{E}_{X'}[\hat{H}_n(Y \mid X = X')]$$

$$= \lim_{n \to \infty} \int_{x \in \mathbb{R}^d} \hat{H}_n(Y \mid X = x) \, dF_X,$$

also by the Weak Law. Because of the finiteness of [K], the entropy-like term  $|\hat{H}_n(Y \mid X = x)|$  is bounded by  $\log K$  for all n and x. Therefore, by the Dominated Convergence Theorem,

$$\lim_{n \to \infty} \int_{x \in \mathbb{R}^d} \hat{H}_n(Y \mid X = x) dF_X = \int_{x \in \mathbb{R}^d} \lim_{n \to \infty} \hat{H}_n(Y \mid X = x) dF_X$$
$$= \int_{x \in \mathbb{R}^d} H(Y \mid X = x) dF_X$$
$$= \mathbb{E}_{X'}[H(Y \mid X = X')]$$
$$= H(Y \mid X)$$

Thus,  $\hat{H}(Y \mid X)$  is consistent for  $H(Y \mid X)$ .

Letting the n and  $\nu$  be a fixed fraction of the training set size is a special case of both growing to infinity. The second result concerns the mutual information estimate. The entropy H(Y) is estimated in the natural way, as in

$$\hat{H}(Y) = -\sum_{y \in [K]} \hat{p}(y) \log \hat{p}(y)$$

where  $\hat{p}(y) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\{Y_i = y\}$  (and set to zero when appropriate). Consequently, the mutual information is estimated as

$$\hat{I}_{n,\nu}(X,Y) = \hat{H}(Y) - \hat{H}(Y \mid X).$$

Consider the same assumptions as Theorem 1.

Theorem 2 (Consistency of the mutual information estimate).  $\hat{I}_{n,\nu}(X,Y) \to_p I(X;Y)$  as  $n,\nu \to \infty$ .

*Proof.* For i.i.d. observations of  $Y_i$ , it is clear that  $\hat{H}(Y)$  is a consistent estimate for H(Y). By Theorem 1 we have the consistency of  $\hat{H}(Y\mid X)$  for  $H(Y\mid X)$ . Thus, the consistency of  $\hat{I}(X,Y)$  follows immediately.