# MANIFOLD OBLIQUE RANDOM FORESTS: TOWARDS CLOSING THE GAP ON CONVOLUTIONAL DEEP NETWORKS

ADAM LI*,1,3, RONAN PERRY*,1, CHESTER HUYNH*,1,3, TYLER M. TOMITA1, RONAK MEHTA2, JESÚS ARROYO2, JESSE PATSOLIC2, BENJAMIN FALK2, SRIDEVI V. SARMA1,3, JOSHUA T. VOGELSTEIN1,2,3*

**Abstract.**
Decision forests, in particular random forests and gradient boosting trees have demonstrated state-of-the-art accuracy compared to other methods in many supervised learning scenarios. Forests dominate other methods in tabular data, that is, when the feature space is unstructured, so that the signal is invariant to a permutation of the feature indices. However, in structured data lying on a manifold—such as images, and time-series—deep networks, specifically convolutional deep networks (ConvNets), tend to outperform forests. We conjecture that it is in part due to networks not simply analyzing feature magnitudes, but also their indices. In contrast, naïve forest implementations fail to explicitly consider feature indices. A recent approach demonstrates that forests, for each node, implicitly sample a random matrix from some specific distribution. These forests, like some networks, learn by partitioning the feature space into convex polytopes corresponding to linear functions. We build on that approach with Manifold Oblique Random Forests (Morf) that chooses distributions in a *manifold-aware fashion* to incorporate feature locality. Morf runs fast and maintains interpretability and theoretical justification. Morf also has excellent empirical classification performance on simulated data and real images and multivariate time-series. It outperforms non-neural network approaches that ignore feature space structure and challenges the performance of ConvNets in some cases.

**1. Introduction.** Decision forests, including random forests and gradient boosting trees, have solidified themselves in the past couple decades as a powerful ensemble learning method in supervised settings [1, 2], including both classification and regression [3]. In classification, each forest is a collection of decision trees whose individual classifications of a data point are aggregated together using majority vote. One of the strengths of this approach is that each decision tree need only perform better than chance for the forest to be a strong learner, given a few assumptions [4, 5]. Additionally, decision trees are relatively interpretable because they can provide an understanding of which features are most important for correct classification [6]. In 2001, Breiman originally proposed decision trees that partition the data set using hyperplanes aligned to feature axes [6]. Yet, this limits the flexibility of the forest and requires trees of large depth to classify some data sets, leading to overfitting. He also suggested that algorithms which partition based on linear combinations of the coordinate axes can improve performance [6, 7], which was corroborated in subsequent work [8]. More recently, Sparse Projection Oblique Randomer Forest (Sporf)—which leverages sparse random projections of the data—has shown impressive improvement over other methods [9]. Other extensions have led to neural decision forests [10, 11] which attempt to combine the strengths of neural networks and random forests by using differentiable functions at split nodes and leaves, leading to trees which can be learned via backpropagation. Under certain function choices, once learned, these forests turn out to be equivalent to neural networks with many zeroed weights [11].

Random forests and other machine learning algorithms typically operate in a tabular setting, viewing an observation $\boldsymbol{x} = (x_1, \ldots, x_p)^T \in \mathbb{R}^p$ as an unstructured feature vector. In doing so, they neglect the feature indices in settings where the in-

---
* 1 Department of Biomedical Engineering, 2 Center for Imaging Science, 3 Institute for Computational Medicine, Kavli Neuroscience Discovery Institute, Johns Hopkins University, * Indicates co-first authorship with equal contribution and any ordering of these authors is allowed

dices encode additional information. For structured data, e.g. images or time series, traditional decision forests do not incorporate the known local structure. For decision forests to utilize known local structure in data, new features encoding this information must be manually constructed or new splitting criterion must be implemented. Prior research has extended random forests to a variety of computer vision tasks [12, 13, 14, 15] and augmented random forests with structured pixel label information [16]. The decision tree at the heart of the Microsoft Kinect showed great success by specializing for image data with depth information [15]. Yet these methods either generate features *a priori* from individual pixels (and thus do not take full advantage of the local topology) or lack the flexibility to learn relevant structure. Other approaches have circumvented the problem of learning from raw structured data through tabular feature engineering, notably employed by the aforementioned deep neural decision forest [10] using a convolutional deep network (ConvNets). Decision forests have also been used to learn distance metrics on unknown manifold structures [17], but such manifold forest algorithms are unsupervised.

Inspired by Sporf, we propose a classification algorithm, Manifold Oblique Random Forests (Morf). Morf takes a projection distribution that accounts for neighboring features on a manifold, while incorporating enough randomness to learn the relevant projections. At each node in the decision tree, a set of neighboring features are randomly selected using knowledge of the underlying manifold. Weighting and summing the values of the selected features yields a set of oblique projections of the data which can then be evaluated to partition the observations. We show Morf's effectiveness across simulated and real-data settings as compared to common classification algorithms. In each case, Morf performs better than non-ConvNet algorithms that lack local feature information, while approaching, or even improving upon ConvNet performance in certain real data applications. Furthermore, the optimized and parallelizable open source implementation of Morf in Python is available at https://neurodata.io/code/.

**2. Background and Related Work.** We will first define the notation and classification framework needed to describe Morf.

**2.1. Classification.** Let $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ be a random sample from the joint distribution $F_{XY}$ and $D_n := \{(x_i, y_i)\}_{i=1}^n$ be our $n$ observed data points where all $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ are drawn from $F_{XY}$. Denote $\mathcal{X} \subseteq \mathbb{R}^p$ as the space of data vectors, and $\mathcal{Y} = \{1, \dots, K\}$ as the space of K class labels. A classifier is a function that assigns to any unseen data point, $X$, a class label $y \in \mathcal{Y}$. Our goal is to learn a classifier $g_n(X; D_n) : \mathcal{X} \times (\mathcal{X} \times \mathcal{Y})^n \to \mathcal{Y}$ from our data $D_n$ that minimizes the expected risk corresponding to 0-1 loss, equivalently the probability of incorrect classification,

$$L(g) := \mathbb{E}[\mathbb{I}[g(X) \neq Y]] = P(g(X) \neq Y),$$

with respect to the distribution of $F_{XY}$. The optimal such classifier is the Bayes classifier

$$g^*(X) := \operatorname*{argmax}_{y \in \{1, \dots, K\}} P(Y = y \mid X),$$

which has the lowest attainable risk $L^* := L(g^*(X))$.

**2.2. Random Forests.** Originally popularized by Amit and Geman [18] and subsequently codified by Breiman, the random forest (RF) classifier is empirically very effective [1] while maintaining strong theoretical guarantees [6]. A random forest is an ensemble of decision trees whose individual classifications of a data point are aggregated together using majority vote. Each decision tree recursively partitions the feature space and then makes separate predictions in each of the final subspaces. A partition occurs at a split node in the tree on a subset of the data $S = \{(x_i, y_i)\} \subseteq D_n$. The node is split into two child nodes, each associated with a partition of $S$ based on the value of a selected feature $j \in \{1, \cdots, p\}$. Let $e_j \in \mathbb{R}^p$ denote a unit vector in the standard basis (that is, a vector with a single one in the $j$th entry and the rest of the entries are zero) and $\tau \in \mathbb{R}$ a threshold value. Then, $S$ is partitioned into the two subsets, a left node $(L)$, and a right node $(R)$.

$$S_\theta^L = \{(x_i, y_i) \mid e_j^\mathsf{T} x_i < \tau, (x_i, y_i) \in S\},$$
$$S_\theta^R = \{(x_i, y_i) \mid e_j^\mathsf{T} x_i \geq \tau, (x_i, y_i) \in S\}$$

given the parameter pair $\theta = \{e_j, \tau\}$. To choose the partition, $\theta$ is sampled $d$ times (also known as $m_{try}$ in the literature). Then the locally optimal $\theta^* = (e_j^*, \tau^*)$ pair is greedily selected from among a set of $d$ randomly selected standard basis vectors as that which maximizes some measure of information gain. A typical measure is a decrease in impurity, calculated by the Gini impurity score $I(S)$, of the resulting partitions [3]. Let $\hat{p}_k(S) = \frac{1}{|S|} \sum_{y_i \in S} \mathbb{I}[y_i = k]$ be the fraction of elements of class $k$ in partition $S$ and $I(S) := \sum_{k=1}^K \hat{p}_k(1 - \hat{p}_k)$ be the Gini impurity. Then the split

$$\theta^* = \operatorname*{argmax}_\theta |S|I(S) - |S_\theta^L|I(S_\theta^L) - |S_\theta^R|I(S_\theta^R)$$

is chosen to maximize the decrease in impurity from the parent node containing $S$. A leaf node in the decision tree is created once a partition reaches a stopping criterion, typically either falling below an impurity score threshold or a minimum number of samples [3].

To classify a feature vector $x$, it is evaluated at root node of the tree and split into one of the two partitions. This process is repeated recursively at subsequent split nodes until $x$ "falls into" a leaf, upon which posterior probability estimates of the class labels can be assigned. Let $l_b(x)$ be the set of training examples at the leaf node in tree $b$ into which $x$ falls. The empirical posterior probability of label $y$ in $b$ is thus $\hat{p}_{n,b}(y \mid x) = \frac{1}{|l_b(x)|} \sum_{i=1}^n \mathbb{I}[y_i = y]\mathbb{I}[x_i \in l_b(x)]$. The forest composed of $B$ trees computes the empirical posterior probability for $x$ by averaging over the trees $\hat{p}_n(y \mid x) = \frac{1}{B} \sum_{b=1}^B p_{n,b}(y \mid x)$ and classifies $x$ per the label with the greatest empirical posterior probability [3]

$$g_n(x) = \operatorname*{argmax}_{y \in \{1, \ldots, K\}} \hat{p}_n(y|x).$$

For good performance of the ensemble, the individual decision trees must be relatively uncorrelated from one another. This is typically done by considering a random subset of features at each split node and training each tree on a bootstrapped subsample of the full training data. Applying these techniques reduces the amount random forests overfit and lowers the upper bound of the generalization error [6].

**2.3. Oblique Forests.** Sparse Projection Oblique Randomer Forests (SPORF), is a recent modification to random forest that has shown improvement over axis-aligned random forests and other oblique forests that compute linear combinations of features [7, 8, 9]. Recall that RF split nodes partition data along the coordinate axes by comparing the projection $e_j^\mathsf{T} x$ of observation $x$ on standard basis $e_j$ to a threshold value $\tau$. SPORF generalizes the set of possible projections, allowing for the data to be partitioned along any linear combination of axes specified by the sparse vector $a_j \in \mathbb{R}^p$. The partition

$$
\begin{aligned}
S_\theta^L &= \{(x_i, y_i) \mid a_j^\mathsf{T} x_i < \tau, (x_i, y_i) \in S\}, \\
S_\theta^R &= \{(x_i, y_i) \mid a_j^\mathsf{T} x_i \geq \tau, (x_i, y_i) \in S\}
\end{aligned}
$$

follows from our choice of $\theta = \{a_j, \tau\}$, where the entries of $a_j$ vector entries are defined as follows (here $a_{ij}$ is the ith entry of $a_j$):

$$
a_{ij} = \begin{cases} 1 & \text{with prob. } \frac{1}{2s} \\ 0 & \text{with prob. } 1 - \frac{1}{s} \\ -1 & \text{with prob. } \frac{1}{2s} \end{cases}
$$

Then $\theta$ is chosen in the same manner as in axis-aligned Random Forests. All other aspects of SPORF are the same as RF.

## 3. Methods.

**3.1. Sampling Projections from a Dictionary.** To move towards manifold forests, we observe that random and oblique forests both sample atoms from a dictionary to create their projected feature values, which are compared with threshold $\tau$ to determine the partition. In axis-aligned random forests, the dictionary, $\mathcal{A} = \{e_j\}_{j=1}^p$ is the set of points along the p-dimensional hypercube, i.e. standard basis vectors in $\mathbb{R}^p$. Then at every node, atoms $e_i \in \mathcal{A}$ from the dictionary are sampled $d$ times.

Similarly, in oblique forests, let the dictionary $\mathcal{A}$ be the set of vectors (atoms) $\{a_j\}$, each atom a p-dimensional vector defining a possible projection $a_j^\mathsf{T} x$. In SPORF, the dictionary $\mathcal{A}$ can be much larger then that of Random Forests, because it includes, for example, all 2-sparse vectors. At each split node, SPORF samples $d$ atoms from $\mathcal{A}$ according to a specified distribution. By default, each of the $d$ atoms are randomly generated with the number of non-zero elements drawn from a Poisson distribution with a specified rate $s$. Then, each of the non-zero elements are uniformly randomly assigned either $+1$ or $-1$. Note that although the size of the dictionary for SPORF is $3^p$ (because each of the $p$ elements could be $-1$, 0, or $+1$), the atoms are sampled from a distribution heavily skewed towards sparsity controlled by the $s$ term.

**3.2. Random Projection Forests on Manifolds.** In the structured data setting, the dictionary of atoms $\mathcal{A} = \{a_j\}$ is modified to take advantage of *a priori* knowledge of feature locality on the underlying manifold on which the data lie. We call this 'Manifold Oblique Random Forest' (MORF). This modification constrains the space of random projection decision trees which can be learned in order to better suit certain classification tasks where relations between features may add information. By constructing features in this way, MORF learns low-level features in the structured data, such as corners in images or spikes in time-series.

As in SPORF, let $\mathcal{A}$ be a dictionary of $m$, $p$-dimensional atoms with probability density or mass function $f_{\mathcal{A}}$ over the $m$ atoms. Each atom $a_j \in \mathcal{A}$ projects an observation $x_i$ to a real number $a_j^T x_i$, where nonzero elements of $a_j$ effectively weight and sum features. At each node in the decision tree, MORF selects the best split according to the Gini index over each candidate atom and threshold pair. It has the same partition functions, $S_\theta^L$ and $S_\theta^R$ as SPORF presented in 2.3. What changes when going from SPORF to MORF? We present two generalizations, selection of the non-zero indices and weights, which allow one to specify any type of manifold structure they want.

*Selection of Nonzero Indices.* While the nonzero indices of each atom in SPORF were mutually independent of one another, the key aspect of MORF lies in the user-specified restriction of possible atoms to take advantage of feature locality (e.g. non-zero indices are dependent). Relations between data features, as in feature locality, may be abstractly encoded as a sampling graph in which each feature represents a node and an edge between two features (an adjacency) permits an atom in $\mathcal{A}$ with nonzero weights on both of those features. In Figure 1, if no two features are adjacent (none related), we recover RF(Fig 1a). If all features are pairwise adjacent (all combinations are possible), we recover SPORF(Fig 1b). At a high-level, MORF introduces user-defined networks in between those of RF and SPORF allowing for some relations but not others (Fig 1c).
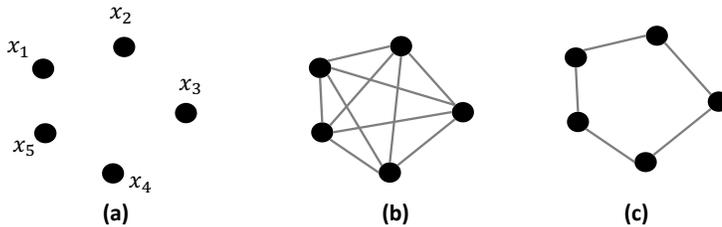


(a)      (b)      (c)

Fig. 1: Intuition of manifold sampling for non-zero indices shown with a 5-dimensional data sample, $x = [x_1, x_2, x_3, x_4, x_5]^T \in \mathbb{R}^5$. Sampling non-zero indices can be abstractly represented as a graph. In axis-aligned random forests (a), at every node, samples from the dictionary are drawn where only one non-zero weight is drawn (i.e. the standard basis vector). This is represented by a completely disconnected graph;. In SPORF (b) any combination of non-zero indices is possible, indicated by a fully-connected graph; each sample can effectively form a linear combination of potentially all features. In practice, the sampled projection vector is sparse, which can be represented by the edge weights being extremely small ($\frac{1}{2s}$). In MORF (c) prior information about the structure of $x$ can be leveraged to constrain the dictionary of possible projection vectors. For example, in natural images, we expect adjacent pixels to be correlated and $x_1, x_2, x_5$ may represent adjacent pixels in a vectorized image. Thus, we would sample non-zero indices from these "patches" of the image.

*Selection of Nonzero Weights.* In RF, all non-zero weights have value 1. In SPORF, all non-zero weights have value 1, or -1. In our settings MORF, all non-zero weights will have value 1. For the setting of natural images, this is equivalent to taking the summation operator of a small "patch". These weights can be set in a more general fashion. For example, one can set weights according to a Gaussian kernel, where

values are higher in the center of the patch and lower towards the edge of the patch. We do not consider these cases, but discuss their implications in Section 7.

**3.2.1. Examples of Projection Dictionaries for Manifolds.** For our applications, we provide a concrete implementation targeted at translation equivariant feature locality in 1D and 2D, such as time series and images respectively. Based on assumptions about the data manifold, one can specify the graph to sample non-zero indices. Assume there exists a similarity matrix, $S$, induced by the manifold dictating how features are related to each other in a data sample, $x_i$. For example, in natural images, $S$ would be a matrix in $\mathbb{R}^{W \times H}$ with width, W, and length, H. So randomly sample $i \in \{1, ..., W\}$ and $j \in \{1, ..., H\}$. Then $S_{ij}$ is the ijth entry of the similarity matrix, corresponding to the ijth feature in $x_i$. We would then randomly sample k-hop neighbors of $S_{ij}$ to form a "patch", and then combine these features using the summation operator (these features can be combined with different weights, which is considered in the Discussion). This patch vectorized would form our projection vector, $a_i$. $\langle x_i^T, a_i \rangle$ would give us a candidate feature value to split on. In practice we do not actually know the similarity matrix induced by the manifold, but we can leverage prior information. For example, in natural images, nearby pixels One can specify the possible widths and heights of sampled patches based on knowledge of the data manifold.

In natural images, at each split node, a set of atoms are randomly sampled to produce candidate features across observations. MORF accepts hyperparameters defining the minimum and maximum number of patch rows $\{h_{min}, h_{max}\}$ and columns $\{w_{min}, w_{max}\}$, respectively. To sample a patch, first the number of rows $h$ and columns $w$ are independently and uniformly sampled between respective minima and maxima (inclusive). As columns are contiguous, a reference leftmost column in the unraveled matrix is sampled as $u \sim \mathcal{U}\{-w + 1, W\}$. If 2D locality is specified, then a reference upper row is sampled as $v \sim \mathcal{U}\{-h + 1, H\}$. In both cases, the reference column and row may be outside of the matrix so that each feature has an equally likely chance of being included in a patch. The region outside of the matrix is ignored, effectively a zero-padded boundary. The algorithm pseudocode is equivalent to that of SPORF except for the distribution $f_A$ described above which can be seen in Appendix C.

In multivariate time-series, the features of a single observation $x_i \in \mathbb{R}^p$ can also be viewed as organized into a 2D matrix in $\mathbb{R}^{S \times T}$, where $S * T = p$. Each feature is indexed by a row and column over the number of sensors (S) and time points (T). In this case, only 1D locality might be beneficial where contiguous time points are correlated. However, unless the sensors are ordered in a meaningful manner, there is no reason to suspect locality along the sensor axis (i.e. columns). Again the parameters of the patch can be specified according to the data manifold.

In our experiments, the atom weights were limited to values of 1 and 0 to limit combinatorial complexity but domain-specific atom design may be desired in some settings along with further task-specific atoms. For graph-valued data, one may consider sampling a collection of neighboring edges or nodes [19]. For spatial related data, one can consider sampling a colleciton of nearby points in Euclidean or Riemannian space [20].

**3.3. Feature Importance.** One of the benefits to decision trees is that their results are fairly interpretable in that they allow for estimation of the relative importance of each feature. Many approaches have been suggested [6, 21], and here we introduce a projection forest specific metric which counts the the number of times

a given feature was used in projections across the ensemble. Formally, a decision tree $T$ in the trained forest $\mathtt{F}$ contains a set of split nodes, where each node $s \in T$ is associated with an atom $a_s^*$ from the dictionary of atoms in the forest $\mathcal{A}_{\mathtt{F}}$ and a threshold that partition the feature space according to the projection $a_s^{*T}x$. Thus, the indices corresponding to nonzero elements of $a_s^*$ indicate important features used in the projection. The importance of feature $k$, denoted $\pi_k$, is calculated as

$$\pi_k = \frac{1}{|\mathcal{A}_{\mathtt{F}}|} \sum_T \sum_{s \in T_S} \mathbb{I}(a_{sk}^* \neq 0),$$

the number of times it is used in a projection, across all decision trees and split. These counts represent the relative importance of each feature in making a correct classification. Such a method applies to $\mathtt{RF}$, $\mathtt{SPORF}$, and $\mathtt{MORF}$ although different results between them would be expected due to different dictionary distributions.

**4. Theoretical Results.** Random forest algorithms have been historically difficult to analyze theoretically, both from a statistical perspectives as well as algorithmic. However, there is a large body of literature making assumptions and modifications on top of Breiman's random forest algorithm [6] from which theoretical analyses are tractable [5, 22, 23, 24, 25, 26, 27, 28]. Here, we provide insights on oblique forests, such as $\mathtt{SPORF}$ and $\mathtt{MORF}$ by expanding upon the axis-aligned forest statistical results in Athey et al. [29] and algorithmic results in Louppe [30].

**4.1. Classifier Consistency.** Athey et al. [29] present a seminal paper specifying some minor distributional assumptions and algorithmic conditions for their generalized random forest algorithm to provide a consistent estimate $\hat{\theta}(x)$ of $\theta(x)$, where $\theta(x)$ is defined as the solution to some estimating equation $\mathbb{E}[\psi_{\theta(x)}(Y_i)|X_i = x] = 0$ for all $x \in \mathcal{X}$ and $\psi_{\theta(x)}$ is a score function. A consistent estimate converges to the true estimand in probability, as the sample size $n$ approaches infinity.

Fundamentally, the consistency of the generalized random forest comes from each tree partitioning the feature space into a set of hyper-rectangles (a bijective map with the set of tree leaves) whose radii go to zero as the sample size grows but slow enough such that they are populated with sets of sizes approaching infinity. This is the same logic behind the consistency of the k-nearest neighbors classifier [31] and indeed random forests are effectively adaptive nearest neighbor classifiers [26].

Although oblique decision trees do not create hyper-rectangular partitions, they do partition the feature space into a finite number of (possibly unbounded) convex polytopes (see Appendix D). Each polytope region responds to a leaf node with a constant classification label per leaf. That oblique trees yield convex polytopes which are not necessarily hyper-rectangles is the only difference compared to axis-aligned trees, and the basis as to why the consistency results of Athey et al. [29] can be extended to the oblique setting.

We show that a main result of Athey et al. [29] holds for oblique random forests, such as $\mathtt{SPORF}$ and $\mathtt{MORF}$. As a corollary, posterior probability estimates are consistent. Thus, an oblique random forest under appropriate conditions admits a consistent classification rule and so its error converges to the minimum expected (Bayes) error (see Appendix A for full proofs). We repeat Specification 1 made by Athey et al. [29] below for reference. Specification 2 is new and restricts the set of possible dictionaries. The full proofs are detailed in Appendix A along with the technical and minor distributional Assumptions 1A-6A.

*Specifications.*
1. All trees are symmetric, in that their output is invariant to permuting the indices of training examples; make balanced splits, in the sense that every split puts at least a fraction $w$ of the observations in the parent node into each child, for some $w > 0$; and are randomized in such a way that, at every split, the probability that the tree splits on the j-th feature is bounded from below by some $\alpha > 0$. The forest is honest and built via subsampling with subsample size $s$ satisfying $s/n \to 0$ and $s \to \infty$ [29].
2. The oblique dictionary $\mathcal{A}$ is finite and contains the set of standard basis vectors $\{e_i\}_{i=1}^{p}$, each with a fixed nonzero probability of being selected at each split node.

*Assumptions.*
1. There exists a density $f$ over $\mathcal{X}$ that is bounded away from zero and infinity. That is, for all $x \in \mathcal{X}$ there exists a $\varepsilon > 0$ such that $\varepsilon < f(x) < \frac{1}{\varepsilon}$.
2. For all $y \in \mathcal{Y}$, $P(Y = y \mid X = x)$ is Lipschitz continuous in $x \in \mathcal{X}$.

*Honesty*, introduced in Specification 1, is a mild condition that removes bias from the leaf estimates by requiring the set of training examples used to learn the structure of the tree to be independent of the set of examples used at the leaf nodes for estimation [5, 22, 24, 29]. In practice this is done using a holdout set per tree and can be beneficial in some cases [28]. Alternatively, this sample splitting can be performed more naturally using the out-of-bag samples from bootstrapping. With the addition of Specification 2, the following theorem extends the results of Athey et al. [29] to oblique regression forests in addition to axis-aligned forests.

THEOREM 1. *Under Assumption 1 and Assumptions 1A-6A [29], the estimate $\hat{\theta}(x)$ from the generalized random forest of Athey et al. [29] incorporating oblique splits and built to Specifications 1-2 is consistent, i.e. $\hat{\theta}_n(x) \xrightarrow{P} \theta(x)$ as $n \to \infty$.*

The conditional mean $\theta(x) = \mathbb{E}[Y_i | X_i = x]$ is a valid estimand for the generalized random forest algorithm and the empirical estimator coincides with that of Breiman's regression forest [29]. So, in the classification setting one may readily estimate the class-conditional posterior in terms of a conditional mean $\theta(x, y) = \mathbb{E}[\mathbb{I}[Y_i = y]|X_i = x] = P(Y_i = y|X_i = x)$ for all $y \in \mathcal{Y}$ in each leaf node. We note that this choice of estimand satisfies Assumptions 2A-6A, see Appendix A for details, with only Assumption 1A remaining as a true assumption. Thus, from Theorem 1 we obtain the following corollary and classification Theorem with Assumption 1 incorporated explicitly.

COROLLARY 2. *Under Assumptions 1-2, posterior estimates from an oblique classification random forest built to Specifications 1-2 are consistent, i.e. $\hat{p}_n(x; y) \xrightarrow{P} P(Y_i = y|X_i = x)$ as $n \to \infty$.*

THEOREM 3. *Under Assumptions 1-2, the classification rule from a oblique classification random forest built to Specifications 1-2 is consistent, i.e. $L_n \xrightarrow{P} L^*$ as $n \to \infty$.*

Note that these results apply to both oblique forests with unstructured atoms such as SPORF as well as those with structured atoms such as MORF. This Lipschitz assumption is a frequent one taken in the literature on random forests [29]. It intuitively makes sense *a priori* that small deviations in $x$ should lead to small deviations in the class probability. As in Athey et al. [29], however, this theorem is limited to continuous-valued features which rules out certain classes of data.

**4.2. Training Time Complexity.** Theoretical analyses are difficult without making assumptions on the data as a tree's possible structure occupies a combinatorially large space and the worst case is too large to be helpful in typical scenarios. We extend the work of Louppe [30] and examine a simplified setting in which the possible sizes induced at each partition node are equally likely. It has been posited and supported empirically that this is a lower bound for the true average case in a RF [30]. One reason that worse-than-average cases may occur is that when none of the candidate features are informative, edge splits are frequent and lead to deep trees [30]. The candidate features in SPORF and MORF are combinations of individual features and we expect this greater flexibility to reduce the chance that no candidate features are informative.

Let a forest have $T$ trees, $d$ candidate features at each split node, and $n$ training samples. At each split node in RF, the complexity is $O(dn \log n)$ to sort observations along each feature using an optimal sorting algorithm [30]. In the average case, the time complexity for RF is then $\mathcal{O}(Tdn \log^2 n)$ [30]. MORF, like SPORF, utilizes sparse matrix multiplication to compute weighted sum while sorting observations at each split node. Thus, letting $H$ and $W$ denote the maximum height and width of a patch in MORF, the time complexity for MORF is $\mathcal{O}(TdHWn \log^2 n)$, because for each of the $d$ features we must make $HW$ multiplication operations. However, as we will show empirically, MORF can find better partitions and thus learn smaller trees.

It is worth noting that because our analysis is based on the framework devised by [29]. The assumptions are rather strict and it is not necessary in all cases for example that the joint density of the explanatory variables is uniformly bounded away from zero and infinity. Nevertheless, our results show that structured oblique methods, such as MORF fit into the theory posed by Athey et al. Moreover, the theory posed by Athey et al. are the basis of the heavily used Generalized Random Forest method and code. Although the study of MORF 's convergence behavior is beyond the scope of this paper, empirically it seems that MORF is able to learn less complex trees compared to RF.

**5. Simulation Experiments.** We examine the performance of MORF in terms of predictive accuracy and runtime in three simulations highlighting 1D and 2D manifolds. In all cases. MORF outperforms methods that do not consider feature locality as well as ConvNets in some cases.

**5.1. Three simulated manifolds.** We evaluated MORF in three simulation settings to show its ability to take advantage of structure in data. MORF was compared to a set of traditional classifiers (and SPORF) that learn from the raw features. For each experiment, we used our open source implementation of MORF as well as SPORF and the RF implementation contained in the SPORF, each with 500 trees. Other classifiers were run from the Scikit-learn Python package [32] and the gradient boosted tree XGBoost (XGB) was run using its Python implementation [33]. Additionally, we tested against a Convolutional Deep Network (ConvNet) built using PyTorch [34] with two convolution layers, ReLU activations, and maxpooling, followed by dropout and a densely connected hidden layer.

Method hyper-parameters were left as defaults except for the ConvNets and MORF which are each specific to the structure of the data and so must be changed. Thus, a well-performing ConvNet architecture was selected and MORF minimum and maximum patch sizes were optimized using a grid search over a range of potential values as well as well as the number of features considered per split (*mtry*) which should vary with the patch size. See Appendix D for details on the hyperparameters and network

architectures across experiments.

Experiment (A) is a non-Euclidean cyclic manifold example inspired by Younes [35] in which the discriminating information is solely contained in the structure of the data. Each observation is a discretization of a circle's perimeter into a one dimensional feature vector with 100 features and two non-adjacent segments of 1's in two differing patterns: class 1 features two segments of length five, while class 2 features one segment of length four and one of length six. Because features are arranged on a circle, segments can wrap around the cyclic feature vector. Figure 2(A) shows examples from the two classes and classification results across various sample sizes.

Experiment (B) is a simple $28 \times 28$ binary image classification problem. Images in class 0 contain randomly sized and spaced *horizontal* bars while those in class 1 contain randomly sized and spaced *vertical* bars. For each sampled image, $k \sim \text{Poisson}(\lambda = 10)$ bars were distributed among the rows or columns, depending on the class. The distributions of the two classes are identical if a 90 degree rotation is applied to one of the classes and so a classifier cannot be learned without learning the structure of the data. Figure 2(B) shows examples from the two classes and classification results across various sample sizes.

Experiment (C) is a signal classification problem highlighting the presence of structure in time series data. One class consists of 100 values of Gaussian noise independent and identically distributed (iid) while the second class has an added exponentially decaying unit step ($u$) beginning at time 20.

$$
\begin{aligned}
X_t^{(0)} &= \epsilon \\
X_t^{(1)} &= u(t - 20) \exp^{(t-20)} + \epsilon, \qquad \epsilon \overset{iid}{\sim} \mathcal{N}(0, 1)
\end{aligned}
$$

Figure 2(C) shows examples from the two classes and classification results across various sample sizes.

In all three simulation settings, MORF outperforms all other classifiers that ignore the local structure, doing especially better at low sample sizes. As compared with ConvNets, MORF sometimes does better, and sometimes worse. The variance across five repeated runs was negligible across sample sizes. The performance of MORF and ConvNets are particularly good in the discretized circle simulation for which most other classifiers perform at chance levels. MORF dominates in the signal classification problem for all sample sizes, most likely because of the ability to learn wide patches which approaches the Bayes classifier. Results on these experiments using uniformly distributed atom weights in between 0 and 1 showed no improvement but increased training time and so were omitted.

We compare the empirical complexity of MORF, SPORF, and RF in Figure 3. Although projection forests required more computations at each partition node during training, as outlined in Section 4.2, MORF is able to learn less complex trees in all cases and sample sizes.

Each simulated experiment was run on CPUs and allocated 45 cores for parallel processing. The resulting train and test times as a function of the number of training samples are plotted in Figure S1. MORF has train and test times slightly longer than those of SPORF. This cost comes at the benefit of less complex trees, as show in in Figure 3. The other method to utilize feature locality, the ConvNet, took noticeably longer to run across simulations for the majority of sample sizes, as seen in Figure S1. Thus its strong performance in those settings comes at an added computational cost, a typical issue for deep learning methods [36].
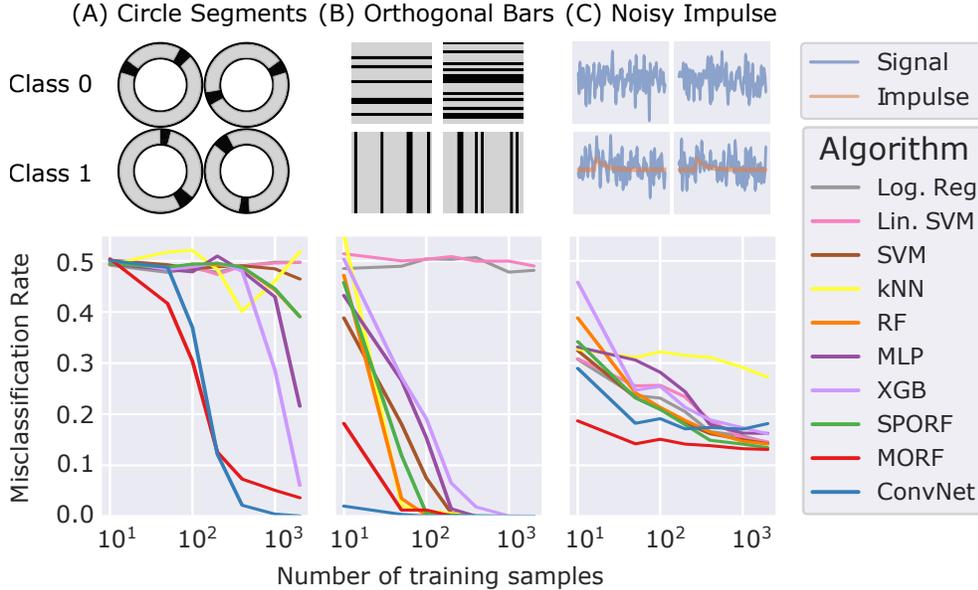
10

Fig. 2: MORF outperforms other algorithms in three two-class classification settings when considering a small number of samples. Upper row shows examples of simulated data from each setting and class. Lower row shows misclassification rate in each setting, tested on 10,000 test samples. **(A)** Two segments in a discretized circle. Segment lengths vary by class. **(B)** Image setting with uniformly distributed horizontal or vertical bars. **(C)** White noise (class 0) vs. exponentially decaying unit impulse plus white noise (class 1). We also observe that ConvNets perform generally better for very large sample sizes (e.g. $>> 10^3$), as expected.

*Model Misspecification.* In the three simulation settings, we explored how MORF is robust to misspecified manifold structure. For example, in (A) Circle Segments, we know that the differentiating factor is the length of the segments being five, or not five. Therefore, one would suspect that the correct patch dimensions should cover that case sufficiently. In Supplementary Figure S3, we see that MORF is relatively stable even when the patch dimensions are not fully the same. The important factor is making sure the information relevant to the task is contained within the possible patch. For full details of the model misspecification experiment, see Supplementary Section B.5.

**5.2. A simulated multivariate time-series.** Experiment (D) here demonstrates how multivariate data with implicit structure can be learned by MORF. The simulation is run as in the prior simulations in Section 5.1. We simulate a multivariate time-series problem, where the signals are governed by a linear dynamical system of the form

$$x(t+1) = Ax(t) + B_i u(t),$$

where the governing linear state matrix, $A \in \mathbb{R}^{3 \times 3}$, is the same, but the class separation is modeled by the input matrices, $B_i \in \mathbb{R}^{3 \times 3}$, defined below. The input to the system, $u(t)$, is the same for all classes. The system is set up such that the pair
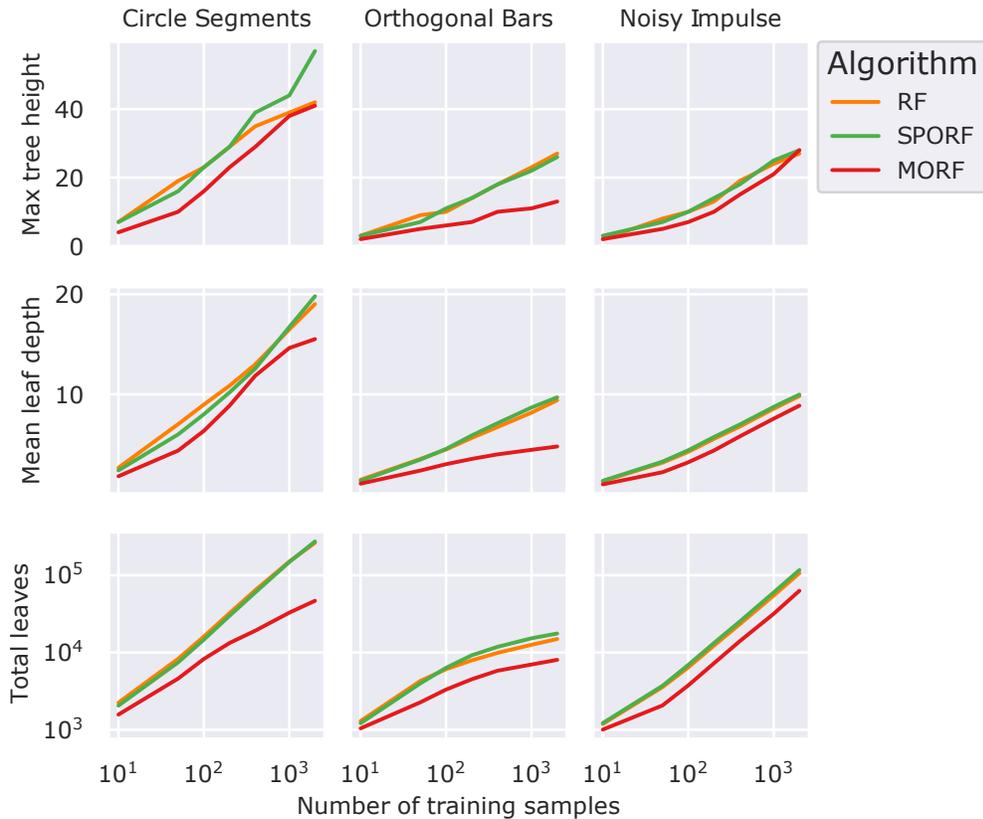
Fig. 3: Empirical random forest complexities in each simulation with respect to their maximum tree height, mean leaf depth, and total number of leaves across training sample size. MORF is able to learn simpler trees in all cases due to its restricted projection distribution.

$(A, B)$ is controllable, and that $A$ is marginally stable (i.e. eigenvalues, $|\lambda| \leq 1$ for all eigenvalues of A.

$$B_0 = \begin{pmatrix} 0.5 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0.5 \end{pmatrix}$$

$$B_1 = \begin{pmatrix} 0.5 & 0.1 & 0.1 \\ 0.1 & 0.5 & 0.1 \\ 0.1 & 0.1 & 0.5 \end{pmatrix}$$

An input $u(t)$ was applied at a random time point selected between the 20th and 40th time point of the simulated time series. A total of 100 time points were simulated and MORF was compared to a suite of other classification algorithms over varying sample sizes. Note that the signals have the exact same dynamics encoded through the $A$ matrix, but the input-output relationships are different. This simulates a setting common in multivariate time-series classification (e.g. EEG, see Section 6.2) where there are signals collected over time that one hypothesizes to be relevant to

12

a task, yet the researcher does not know *a priori* what signals are actually relevant and so they collect as much data as possible. This results in the standard curse-of-dimensionality. However, it is assumed that signals relevant to the task live on a low-dimensional manifold, and the goal is to have a model learn the structure of this manifold for the sake of classification. Figure 4 shows examples from the two classes and classification results across various sample sizes.
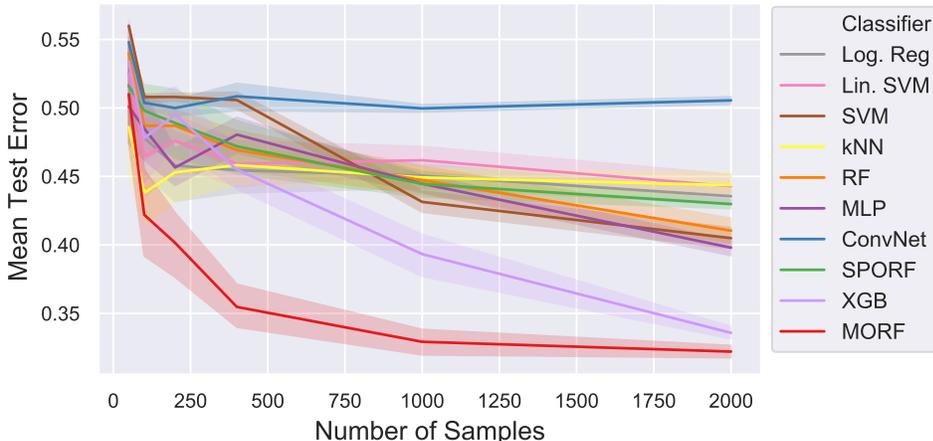


Fig. 4: **Multivariate data embedded in a manifold** Algorithm comparisons on classifying two classes of multivariate data. The data consists of samples from two classes of a 3-dimensional stable linear dynamical system with input. Samples are constructed as $X_i \in \mathbb{R}^{3 \times T}$, where now data points over time are correlated. MORF learns the structure significantly faster then the other classifiers, with XGBOOST requiring more sample to achieve the same test error rate.

A challenging aspect of experiment (D) is that i) the linear state dynamics governed by the A matrix are the same and is considerably larger in norm compared to $B_i$, and ii) the time at which $u(t)$ is applied is random within a small interval. This simulation setting motivates settings where there is a dynamical system with input, such as electroencephalogram (EEG) at rest with dynamics modeled as a linear system, and then a stimulus is applied in the form of input $u(t)$ [37, 38, 39, 40]. Depending on the stimulus applied, this might affect the system in different ways through $B_i$. This is very general setting in where the stimulus can be a flash of light, or indication of movement, or even a direct stimuli to evoke seizures.

**6. Experiments on Real Data.** We next evaluated MORF on three real data sets with varying manifold structure, sample sizes and classification goals. In each dataset, there is the notion of either a 1D or 2D manifold on which we have *a priori* knowledge of feature locality, time and images respectively. We compare results against a suite of classification algorithms as before.

**6.1. 2D Locality: MNIST Digit Classification.** MORF's performance was evaluated on the MNIST dataset, a collection of handwritten digits stored in 28 by 28 square images [41], and compared to the algorithms used in the simulations. 10,000 images were held out for testing and the remaining 50,000 images were used for train-

ing. The results are displayed in Figure 5 (top). Hyperparameters are as described in the simulations, see Appendix D for details. MORF showed an improvement over the other algorithms as compared to ConvNets.
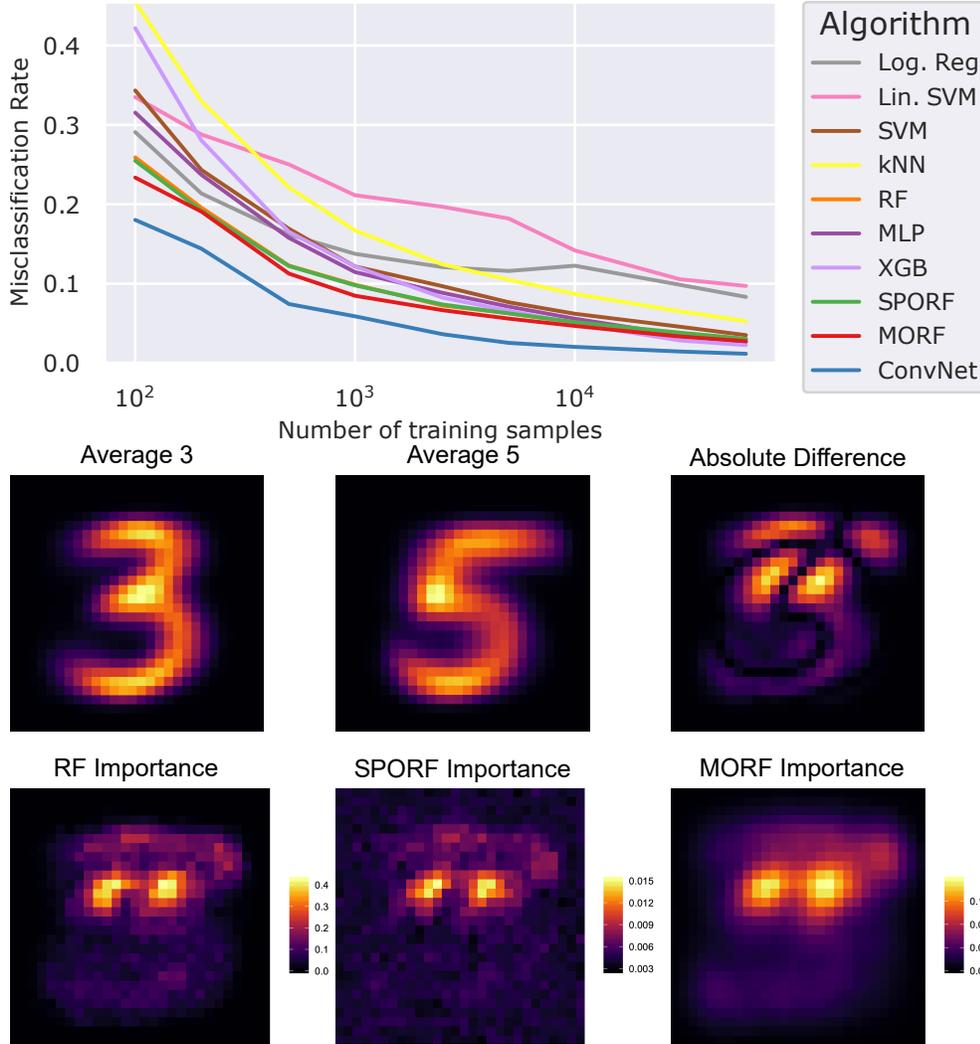


Fig. 5: (**Top**) MORF performance on the MNIST digit classification problem improves prediction accuracy over all other non-ConvNet algorithms, notably in small sample sizes. (**Bottom**) The averages all images from MNIST labeled 3 and 5, respectively, and their absolute difference (top row). Feature importance from MORF (bottom right) shows less noise than SPORF (bottom middle) and is smoother than RF (bottom left).

We then evaluated the ability of MORF to identify important features in manifold-valued data as compared to SPORF and RF. All methods were run on a subset of the MNIST dataset: we only used threes and fives, 100 images from each class.

The feature importance of each pixel is shown in Figure 5 (bottom). MORF visibly results in a smoother pixel importance, a result most likely from the continuity of

neighboring pixels in selected projections. Although Tomita et al. [9] demonstrated empirical improvement of SPORF over RF on the MNIST data, its projection distribution yields scattered importance of unimportant background pixels as compared to RF. Since projections in SPORF have no continuity constraint, those that select high importance pixels will also select pixels of low importance by chance. This may be a nonissue asymptotically, but is a relevant problem in low sample size settings. MORF, however, shows little or no importance of these background pixels by virtue of the modified projection distribution.

**6.2. 1D Locality: Multivariate Time-Series EEG.** MORF's performance was next evaluated on multivariate time-series. Just using raw EEG data, we used MORF to classify movement direction. Compared to a suite of other classification algorithms, MORF is able to achieve a superior performance measured by AUC relative to the other classifiers, as seen in Supplementary Figure S2. We observe that ConvNets have a $0.51 \pm 0.04$ AUC, indicating that it overfit to the training data. This is most likely due to the fact that the dataset presented in [37] consists of 100-200 trials of data (i.e. samples). We next demonstrate that with some structured feature engineering, one can improve MORF on difficult problems beyond that of any traditional classifier and is superior to ConvNets in low-sample size settings.

We next look at a 91 epilepsy subject dataset [42, 43, 44] comprised of intracranial electroencephalogram (iEEG) recordings of multiple seizures. Clinicians annotated a subset of implanted electrodes as part of the clinically hypothesized epileptogenic zone (EZ), and then perform subsequent surgery to resect a super set of those regions. The classification task is to predict surgical outcome of success (seizure free) or failure (seizure recurrence) after a surgical resection is performed on drug resistant epilepsy patients conditioned on the clinically hypothesized EZ regions. If a feature is informative, then it will highly correlate with the clinical EZ when a surgical outcome is successful and vice versa when not successful. In Li et al. [42], Li et al. [43], a feature of the data, "neural fragility" was computed from the data, which is represented as a spatiotemporal heatmap of channels-by-time.

The classification task specifically takes in a data points that are $X_i \in \mathbb{R}^{H \times W}$ of dimension $(20 \times 105)$, where there are 20 quantiles summarizing a distribution of neural fragility and 105 time points around seizure onset. There are 10 quantiles for the neural fragility of electrodes in the clinically hypothesized EZ and 10 quantiles for the neural fragility of electrodes in the rest of the implanted electrodes. The goal is to take $X_i$ and predict $y_i \in \{0, 1\}$, where 0 stands for failed surgical outcome and 1 stands for successful surgical outcome. The classification is setup this way because the clinically annotated EZ electrodes are imperfect and not always representative of the true underlying EZ, which is not observable. Moreover, the multivariate EEG time-series were transformed in this way in order to faciliate comparisons of predictions among subjects with different number of electrodes (ranging from ~30-150). In the original paper, there is preprocessing of the data in the form of a thresholding step, which was done to improve the model performance. However, in this setting, we perform no preprocessing before the classification task to demonstrate fair performance on the "raw" transformed data. For full details on the dataset and clinical problem, we refer the readers to Li et al. [42], Li et al. [43].

In Figure 6a, we compare MORF with default hyperparameters, specified in the SPORF package, against standard classification algorithms as in the simulations. The sample sizes for training are relatively low with 60% of subjects used for training and the rest for the held-out test set. The set of subjects in each cross-validation

index are the same across all classifiers, thus enabling a fair comparison. In some folds, it is seen that all the classifiers perform very poorly. This occurs most likely because the data are noisy, the implanted iEEG electrodes are not perfect and some subjects are very difficult to treat. Moreover, there are only 91 subjects total used in this classification task. Even in this challenging classification setting, we observe that MORF is able to achieve a superior performance measured by AUC (Figure 6b). In terms of the Cohen's effect size, MORF is significantly (p-value $\leq 0.05$) more accurate than all other algorithms besides SPORF per a Wilcoxon paired sign test across the 10 cross-validation folds. We observe that ConvNets perform at chance level on the test set, completely overfitting to the training set in this limited sample size setting.
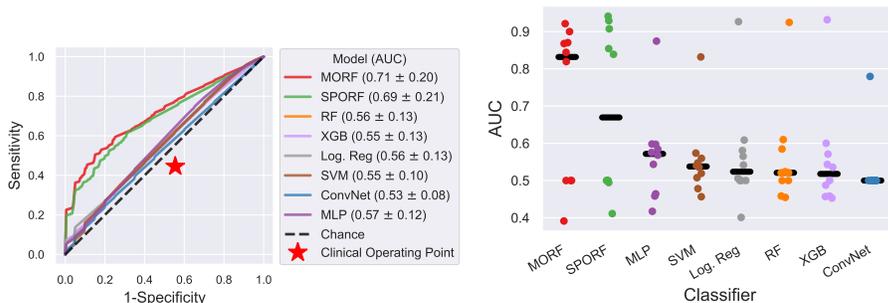


Fig. 6: MORF performance on the epilepsy seizure outcome prediction problem improves prediction accuracy over all other algorithms in 10-fold cross validation (CV). **(Top)** Shows a ROC curve with the mean ROC curve plotted AUC for each classifier. **(Bottom)** Shows a strip plot of the AUC values in 10-fold CV with the median marked in each setting (solid black line). Note that even the ConvNet and multi-layer perceptron (MLP) perform poorly, most likely because the sample size is very low. Compared to the next best non SPORF classifier, MORF improves over the MLP with a Cohen's D effect size of 0.83 (95% CI = [2.32, -0.102]).

**7. Discussion.** The success of sparse oblique projections in decision forests has opened up many possible ways to improve axis-aligned decision forests (including random forests and gradient boosting trees) by way of specialized projection distributions. Traditional decision forests have already been applied to some manifold-valued data, using predefined features to classify images or pixels, and have shown great success, but ignore feature continuity and specialize for specific data modalities. We expand upon sparse oblique projections and introduced manifold-aware projection distributions that exploits prior knowledge of the local topology of a feature space to improve learning rates and accuracy for classification. The open source implementation of MORF subsumes SPORF and provides a flexible classification method for a variety of data modalities and tailored projection dictionaries. We showed in various settings that appropriate domain knowledge can improve the projection distribution and better match ConvNet results (or even outperforming ConvNets significantly) while maintaining interpretability, fast run time, and theoretical justification.

It is plausible that one could design a loss, or engineer a feature based on the structure of one's dataset. Then presumably recent extensions of RF would work even better than MORF on structured data (e.g. data that is spatially dependent) [20]. However, this explicit structure is often not known in practice and hence incorporating

dependence structure directly within the loss function is not as useful for an off-the-shelf tool. Morf circumvents this issue by sampling projection vectors from a dictionary and using a recursive surrogate loss instead (i.e. the Gini impurity per leaf). Moreover, Morf is significantly cheaper in terms of computational cost since we are not directly performing optimization on a desired loss function (e.g. see Supplementary Figure S1 on training and testing times).

The flexibility in choices of Morf's dictionary opens a much larger combinatorial space to sample from compared to a traditional random forest. More complex possibilities may lead to improved performance, but potentially at the cost of greater sampling requirements. Similarly, research into other task-specific projection dictionaries may lead to improved results in computer vision tasks, through better texture quantification for instance, or in other manifold-valued settings such as graphs. Although, unlike ConvNets, Morf is not globally translation equivariant, it can be locally translation equivariant given atoms reminiscent of Gabor filters, for instance. Without local equivariance, discriminative features must be constant in their indices or the training data must be rich enough to fully encapsulate possible observations [15]. The Morf projection distributions may also be incorporated into other state of the art Forest algorithms such as XgBoost. Additionally, the fact that oblique decision forests lead to partitions of convex polytopes is of interest in that it has been shown that deep nets with Rectified Linear Units (ReLUs) or hard tanh activation layers also partition the feature space into convex polytopes with different linear functions on each region [45]. This shared "partition and vote scheme" offers insight into their relationship with one another as well as the functioning of the brain [46]

### References.

[1] M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim, "Do we need hundreds of classifiers to solve real world classification problems?" *Journal of Machine Learning Research*, vol. 15, pp. 3133–3181, 2014.

[2] R. Caruana and A. Niculescu-Mizil, "An empirical comparison of supervised learning algorithms," in *Proc. of the 23rd Int. Conf. on Machine Learning*, ser. ICML '06. New York, NY, USA: ACM, 2006, pp. 161–168.

[3] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, ser. Springer Series in Statistics. New York, NY, USA: Springer New York Inc., 2001.

[4] R. E. Schapire, "The strength of weak learnability," *Mach. Learn.*, vol. 5, no. 2, pp. 197–227, Jul. 1990.

[5] G. Biau, L. Devroye, and G. Lugosi, "Consistency of Random Forests and Other Averaging Classifiers," *Journal of Machine Learning Research*, vol. 9, pp. 2015–2033, Jun. 2008.

[6] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, Oct 2001.

[7] B. H. Menze, B. M. Kelm, D. N. Splitthoff, U. Koethe, and F. A. Hamprecht, "On oblique random forests," in *Machine Learning and Knowledge Discovery in Databases*, D. Gunopulos, T. Hofmann, D. Malerba, and M. Vazirgiannis, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 453–469.

[8] T. M. Tomita, M. Maggioni, and J. T. Vogelstein, "Roflmao: Robust oblique forests with linear matrix operations," in *Proc. of the 2017 SIAM Int. Conf. on Data Mining*, 2017, pp. 498–506.

[9] T. M. Tomita, J. Browne, C. Shen, J. Chung, J. L. Patsolic, B. Falk, C. E. Priebe, J. Yim, R. Burns, M. Maggioni, and J. T. Vogelstein, "Sparse projection oblique

randomer forests," *Journal of Machine Learning Research*, vol. 21, no. 104, pp. 1–39, 2020. [Online]. Available: http://jmlr.org/papers/v21/18-664.html

[10] P. Kontschieder, M. Fiterau, A. Criminisi, and S. R. Bulo, "Deep Neural Decision Forests," in *2015 IEEE Int. Conf. on Computer Vision (ICCV)*. Santiago, Chile: IEEE, Dec. 2015, pp. 1467–1475.

[11] G. Biau, E. Scornet, and J. Welbl, "Neural Random Forests," *arXiv:1604.07143 [cs, math, stat]*, Apr. 2018, arXiv: 1604.07143.

[12] V. Lepetit, P. Lagger, and P. Fua, "Randomized trees for real-time keypoint recognition," in *2005 IEEE Computer Society Conf. on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 2, June 2005, pp. 775–781 vol. 2.

[13] J. Gall, A. Yao, N. Razavi, L. Van Gool, and V. Lempitsky, "Hough forests for object detection, tracking, and action recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 11, pp. 2188–2202, Nov 2011.

[14] A. Bosch, A. Zisserman, and X. Munoz, "Image classification using random forests and ferns," in *2007 IEEE 11th Int. Conf. on Computer Vision*, Oct 2007, pp. 1–8.

[15] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," in *CVPR 2011*, June 2011, pp. 1297–1304.

[16] P. Kontschieder, S. R. Bulò, H. Bischof, and M. Pelillo, "Structured class-labels in random forests for semantic image labelling," in *2011 Int. Conf. on Computer Vision*, Nov 2011, pp. 2190–2197.

[17] A. Criminisi, J. Shotton, and E. Konukoglu, "Decision forests: A unified framework for classification, regression, density estimation, manifold learning and semi-supervised learning," *Found. Trends. Comput. Graph. Vis.*, vol. 7, no. 2&#8211;3, pp. 81–227, Feb. 2012.

[18] Y. Amit and D. Geman, "Shape quantization and recognition with randomized trees," *Neural Comput.*, vol. 9, no. 7, pp. 1545–1588, Oct. 1997.

[19] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, "A comprehensive survey on graph neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, p. 1–21, 2020. [Online]. Available: http://dx.doi.org/10.1109/TNNLS.2020.2978386

[20] A. Saha, S. Basu, and A. Datta, "Random forests for spatially dependent data," *Journal of the American Statistical Association*, vol. 0, no. 0, pp. 1–19, 2021. [Online]. Available: https://doi.org/10.1080/01621459.2021.1950003

[21] S. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *ArXiv*, vol. abs/1705.07874, 2017.

[22] S. Wager and G. Walther, "Adaptive Concentration of Regression Trees, with Application to Random Forests," *arXiv:1503.06388 [math, stat]*, Apr. 2016, arXiv: 1503.06388. [Online]. Available: http://arxiv.org/abs/1503.06388

[23] N. Meinshausen, "Quantile Regression Forests," *Journal of Machine Learning Research*, vol. 7, no. 35, pp. 983–999, 2006. [Online]. Available: http://jmlr.org/papers/v7/meinshausen06a.html

[24] M. Denil, D. Matheson, and N. D. Freitas, "Narrowing the gap: Random forests in theory and in practice," in *Proc. of the 31st Int. Conf. on Machine Learning*, ser. Proc. of Machine Learning Research, E. P. Xing and T. Jebara, Eds., vol. 32. Bejing, China: PMLR, 22–24 Jun 2014, pp. 665–673. [Online]. Available: http://proceedings.mlr.press/v32/denil14.html

[25] G. Biau, "Analysis of a Random Forests Model," *Journal of Machine*

*Learning Research*, vol. 13, no. 38, pp. 1063–1095, 2012. [Online]. Available: http://jmlr.org/papers/v13/biau12a.html

[26] Y. Lin and Y. Jeon, "Random Forests and Adaptive Nearest Neighbors," *Journal of the American Statistical Association*, vol. 101, no. 474, pp. 578–590, 2006, publisher: [American Statistical Association, Taylor & Francis, Ltd.]. [Online]. Available: https://www.jstor.org/stable/27590719

[27] E. Scornet, G. Biau, and J.-P. Vert, "Consistency of random forests," *Ann. Statist.*, vol. 43, no. 4, pp. 1716–1741, Aug. 2015, arXiv: 1405.2881. [Online]. Available: http://arxiv.org/abs/1405.2881

[28] S. Wager and S. Athey, "Estimation and Inference of Heterogeneous Treatment Effects using Random Forests," *Journal of the American Statistical Association*, vol. 113, no. 523, pp. 1228–1242, Jul. 2018. [Online]. Available: https://www.tandfonline.com/doi/full/10.1080/01621459.2017.1319839

[29] S. Athey, J. Tibshirani, and S. Wager, "Generalized random forests," *Ann. Statist.*, vol. 47, no. 2, pp. 1148–1178, 04 2019. [Online]. Available: https://doi.org/10.1214/18-AOS1709

[30] G. Louppe, "Understanding Random Forests: From Theory to Practice," Jun. 2015, arXiv: 1407.7502.

[31] L. Devroye, L. Györfi, and G. Lugosi, *A Probablistic Theory of Pattern Recognition*, 01 1996, vol. 31.

[32] 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*, vol. 12, pp. 2825–2830, 2011.

[33] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proc. of the 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, ser. KDD '16.  New York, NY, USA: Association for Computing Machinery, 2016, p. 785–794. [Online]. Available: https://doi.org/10.1145/2939672.2939785

[34] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.

[35] L. Younes, "Diffeomorphic Learning," *Journal of Machine Learning Research*, vol. 21, no. 220, pp. 1–28, 2020. [Online]. Available: http://jmlr.org/papers/v21/18-415.html

[36] R. Livni, S. Shalev-Shwartz, and O. Shamir, "On the computational efficiency of training neural networks," in *Proc. of the 27th Int. Conf. on Neural Information Processing Systems - Volume 1*, ser. NIPS'14.  Cambridge, MA, USA: MIT Press, 2014, pp. 855–863.

[37] M. S. D. Kerr, P. Sacré, K. Kahn, H.-J. Park, M. Johnson, J. Lee, S. Thompson, J. Bulacio, J. Jones, J. González-Martínez, C. Liégeois-Chauvel, S. V. Sarma, and J. T. Gale, "The Role of Associative Cortices and Hippocampus during Movement Perturbations." *Front. Neural Circuits*, vol. 11, p. 26, 2017. [Online]. Available: http://www.ncbi.nlm.nih.gov/pubmed/28469563http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5395558

[38] A. Li, Z. Fitzgerald, J. Hopp, E. Johnson, N. Crone, J. Bulacio, J. Martinez-Gonzalez, S. Inati, K. Zaghloul, and S. V. Sarma, "Virtual Cortical Stimulation Mapping of Epilepsy Networks to Localize the Epileptogenic Zone," in *Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc. EMBS*, vol. 2019, no. Cc.  Institute of Electrical and Electronics Engineers (IEEE), oct 2019, pp. 2328–2331.

[Online]. Available: https://pubmed.ncbi.nlm.nih.gov/31946366/

[39] A. Li, K. Gunnarsdottir, S. Inati, K. Zaghloul, J. Gale, J. Bulacio, J. Martinez-Gonzalez, and S. Sarma, "Linear time-varying model characterizes invasive EEG signals generated from complex epileptic networks," in *Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc. EMBS*, 2017.

[40] M. Jones, B. McDermott, B. L. Oliveira, A. O'Brien, D. Coogan, M. Lang, N. Moriarty, E. Dowd, L. Quinlan, B. Mc Ginley, E. Dunne, D. Newell, E. Porter, M. A. Elahi, M. O' Halloran, and A. Shahzad, "Gamma Band Light Stimulation in Human Case Studies: Groundwork for Potential Alzheimer's Disease Treatment," *J. Alzheimers. Dis.*, vol. 70, no. 1, pp. 171–185, 2019. [Online]. Available: https://pubmed.ncbi.nlm.nih.gov/31156180https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6700637/

[41] Y. Lecun, C. Cortes, and C. J. Burges, *The MNIST Database of Handwritten Digits*, 1999.

[42] A. Li, C. Huynh, Z. Fitzgerald, I. Cajigas, D. Brusko, J. Jagid, A. Claudio, A. Kanner, J. Hopp, S. Chen, J. Haagensen, E. Johnson, W. Anderson, N. Crone, S. Inati, K. Zaghloul, J. Bulacio, J. Gonzalez-Martinez, and S. V. Sarma, "Neural fragility as an eeg marker of the seizure onset zone," *bioRxiv*, 2021. [Online]. Available: https://www.biorxiv.org/content/early/2021/02/02/862797

[43] A. Li, S. Inati, K. Zaghloul, and S. Sarma, "Fragility in epileptic networks: The epileptogenic zone," in *2017 American Control Conf. (ACC)*, 2017, pp. 2817–2822.

[44] A. Li, P. Myers, N. Warsi, K. M. Gunnarsdottir, S. Kim, V. Jirsa, A. Ochi, H. Otusbo, G. M. Ibrahim, and S. V. Sarma, "Neural fragility of the intracranial eeg network decreases after surgical resection of the epileptogenic zone," *medRxiv*, 2021. [Online]. Available: https://www.medrxiv.org/content/early/2021/07/08/2021.07.07.21259385

[45] M. Raghu, B. Poole, J. Kleinberg, S. Ganguli, and J. Sohl-Dickstein, "On the Expressive Power of Deep Neural Networks," *arXiv:1606.05336 [cs, stat]*, Jun. 2017, arXiv: 1606.05336. [Online]. Available: http://arxiv.org/abs/1606.05336

[46] C. E. Priebe, J. T. Vogelstein, F. Engert, and C. M. White, "Modern Machine Learning: Partition & Vote," *bioRxiv*, p. 2020.04.29.068460, Apr. 2020, publisher: Cold Spring Harbor Laboratory Section: New Results. [Online]. Available: https://www.biorxiv.org/content/10.1101/2020.04.29.068460v1

[47] R. Perry, R. Mehta, R. Guo, J. Arroyo, M. Powell, H. Helm, C. Shen, and J. T. Vogelstein, "Random Forests for Adaptive Nearest Neighbor Estimation of Information-Theoretic Quantities," *arXiv:1907.00325 [cs, stat]*, Sep. 2021, arXiv: 1907.00325. [Online]. Available: http://arxiv.org/abs/1907.00325

[48] M. S. Breault, Z. B. Fitzgerald, P. Sacré, J. T. Gale, S. V. Sarma, and J. A. González-Martínez, "Non-motor brain regions in non-dominant hemisphere are influential in decoding movement speed," *Frontiers in Neuroscience*, vol. 13, p. 715, 2019. [Online]. Available: https://www.frontiersin.org/article/10.3389/fnins.2019.00715

**Appendices.**

**Appendix A. Proofs.**

**A.1. Convex Polytope Partition Results.** As mentioned in Section 4.1, a random projection tree partitions the feature space into a finite number of (possibly unbounded) convex polytopes. The proof of that is as follows.

*Proof.* A convex polytope in $d$ dimensions can be defined as the union of a finite number of halfspaces, where a halfspace is a $d-1$ dimensional surface defined by the linear inequality

$$a^T x \leq b$$

for fixed $a \in \mathbb{R}^d$ and $b \in \mathbb{R}$. In a random projection tree, each split node $i$ partitions the set of points at that node according to such an inequality $a_i^T x \leq b_i$. Consider the path of $k$ split nodes, including the root, to a leaf $l$ and the set of corresponding halfspace defining $\{(a_i, b_i)\}_{i=1}^k$ terms for each split node. We see that in the feature space $S$, the subset that "falls into" leaf $l$ is the solution set to

$$A_l x \leq b_l$$

where $A_l = [a_1, \ldots, a_k]^T$ and $b_l = [b_1, \ldots, b_k]^T$.

Thus each leaf node forms a convex polytope. Additionally, note that any $x \in S$ will deterministically end up in a leaf node (by classification of $x$) as the tree is of finite depth and that all leaf node convex polytopes are mutually exclusive as the lowest common ancestor of any two leaves forms mutually exclusive sets. If the feature space is unbounded, then at least one partition must be unbounded too. Thus, a tree partitions the feature space into a finite number of possibly infinite convex polytopes.□

**A.2. Consistency Results.** The least we can ask of our classification rule $\{g_n\}_{n=1}^\infty$ is for it to be consistent, $L_n \overset{P}{\to} L^*$ as $n \to \infty$, where $L_n$ and $L^*$ are the expected 0-1 losses of the finite sample rule $g_n$ and the Bayes decision rule $g^*$, respectively. Our results here build upon the results of Athey et al. [29] who prove under Assumption 1, 1A-6A and Specification 1 that their generalized random forest (GRF) algorithm, which subsumes Breiman's regression forest, provides a consistent estimate $\hat{\theta}_n(x)$ for some quantity $\theta(x)$. The estimand $\theta(x)$ is defined as the solution to the generic estimating equation $M_\theta(x) := \mathbb{E}[\psi_\theta(Y_i)|X_i = x] = 0$ for all $x \in \mathcal{X}$ where $\psi_\theta$ is a score function. We begin by restating Assumptions 1A-6A of the GRF algorithm as they are relevant to further results but strongly recommend referring to the original paper [29] for additional details.

*Assumptions:.*

1A. For fixed values $\theta(x)$, we assume that $M_\theta(x)$ is Lipschitz continuous in $x$.

2A. When $x$ is fixed, we assume that $M_\theta(x)$ is twice continuously differentiable in $\theta$ with a uniformly bounded second derivative, and that $\frac{\partial}{\partial(\theta)} M_\theta(x) |_\theta \neq 0$ for all $x \in \mathcal{X}$.

3A. The worst-case variogram of $\psi_\theta(Y)$ is Lipschitz-continuous in $\theta(x)$.

4A. The $\psi$-functions can be written as $\psi_\theta(Y) = \lambda(\theta(x); Y) + \xi_\theta(g(Y))$, such that $\lambda$ is Lipschitz-continuous in $\theta$, $g : \{Y\} \to \mathrm{R}$ is a univariate summary of $Y$, and $\xi_\theta : \mathrm{R} \to \mathrm{R}$ is any family of monotone and bounded functions.

5A. For any weights $\alpha_i(x)$ such that $\sum_i \alpha_i(x) = 1$, the estimation equation returned a minimizer $\hat{\theta}(x)$ that at least approximately solves the estimating equation $||\sum_{i=1}^n \alpha_i(x)\psi_{\hat{\theta}}(Y_i)||_2 \leq C max\{\alpha_i(x)\}$ for some constant $C \geq 0$.

6A. The score function $\psi_\theta(Y)$ is a negative sub gradient of a convex function, and the expected score $M_\theta(x)$ is the negative gradient of a strongly convex function.

*Proof of Theorem 1: Consistent oblique random forests posterior estimates.* Our Theorem 1 extends Theorem 3 of the GRF paper [29] given the additional Specification 2. Specifically, the GRF Theorem 3 in Athey et al. [29] result relies on the consistency result of Theorem 1 in Wager and Athey [28]. We need only to verify foundations of Theorem 1 [28] that take into account the use of axis-aligned splits, those being Theorems 3 and 5 of Wager and Athey [28]. Thus, it suffices to confirm that those results are unchanged under oblique splits and given Specification 2 holds.

Theorem 3 [28] proves an asymptotic upper bound on the diameter of a leaf, $\mathrm{diam}(L(x))$, by applying an asymptotic upper bound result from Lemma 2 [28] on the diameter of dimension $j$ in the leaf, $\mathrm{diam}_j(L(x))$. The leaf $L(x)$ is a polytope formed from a combination of axis-aligned and oblique splits. Considering only the axis-aligned conditions forming $L(x)$, by the positive probability of splitting on each dimension per Specification 2, the upper bound of Lemma 2 Wager and Athey [28] holds. As the addition of oblique conditions cannot increase the size of the leaf, the same upper bound holds. Similarly, the diameter $\mathrm{diam}(L(x))$ of the leaf is smaller than the diameter of the polytope formed from just axis-aligned conditions. By the diameter bound from Lemma 2 [22] of each feature, the upper bound of Theorem 3 [28] holds for the axis-aligned polytope and so also $L(x)$.

Theorem 5 [28] hinges on Lemma 4 [28] which brings up the concept of a *potential nearest neighbor* (PNN) [22, 28].

DEFINITION 1. *$x_i \in \{x_1, \ldots, x_s\} \in \{\mathbb{R}^p\}^s$ is a potential nearest neighbor (PNN) of $x \in \mathbb{R}^p$, if there is an axis-aligned hyperrectangle containing only $x$ and $x_i$. A k-PNN set is a collection of $k$ points and $x$ in an axis-aligned hyperrectangle containing no other points. A predictor $T$ for $x$ is a k-PNN predictor if given*

$$\{z\} = \{(x_1, y_1), \ldots, (x_s, y_s)\} \in \{\mathbb{R}^p \times \mathcal{Y}\}^s,$$

*$T$ outputs the average of the $y_i$ among a k-PNN set of $x$ with respect to the $x_i$.*

In the case of oblique split decision trees, we have the following result.

LEMMA 1. *Let $T$ be a decision tree which makes oblique splits (including axis-aligned splits) at each interior node with finite dictionary $\mathcal{A}$ of $m$ vectors encoding the set of allowable oblique axes. If $T$ has leaves between size $k$ and $2k - 1$, then $T$ is a k-PNN predictor on $\mathbb{R}^m$.*

*Proof.* Let $\mathcal{X}$ denote the vector space of possible samples, where $x \in \mathcal{X} \subset \mathbb{R}^p$. Since $\mathcal{A} \in \{\mathbb{R}^p\}^m$, let $A \in \mathbb{R}^{p \times m}$ denote the matrix whose columns are the elements of $\mathcal{A}$. Then $\mathcal{B} = A^T \mathcal{X} \subset \mathbb{R}^m$ is a vector space of dimension at most $\min(p, m)$ in a space of dimension $m$. Bases of $\mathcal{B}$ correspond to bases or oblique combinations of them from $\mathcal{X}$ and so every oblique split in $\mathcal{X}$ is an axis-aligned split in $\mathcal{B}$. The points which fall into a leaf of $T$ are the only points which satisfy the linear system formed

by the set of splits, which are the only points that fall into the hyperrectangle in $\mathcal{B}$ defined by that system. As any decision tree making axis-aligned splits with leaves of sizes between $k$ and $2k-1$ is a $k$-PNN predictor [26], $T$ is thus a $k$-PNN predictor in $\mathcal{B} \subset \mathbb{R}^m$. □

In this expanded feature space from which we can view oblique splits as axis-aligned, as in the above proof, we can scale down the marginals to be within $[0, 1]$. While this space no longer satisfies the Lipschitz criteria and if $m > p$ may have a density of 0 at all points outside of the $p$ dimensional subspace, Lemma 4 [28] requires neither of these conditions from the original assumptions. So it holds, albeit with the finite constant $m$ instead of $p$. Thus Theorem 5 [28] holds with simply a modified constant which doesn't change the final established asymptotics in Theorem 1 [28]. So Theorem 1 [28] holds in our oblique forest setting and we can extend Theorem 3 of Athey et al. [29] to the oblique setting per our Theorem 1.

*Proof of Corollary 2: Consistent posterior probability estimates.* In the spirit of Perry et al. [47], we now prove the corollary of our previous Theorem that the consistent regression estimate results of Athey et al. [29] extend to classification and so oblique random forests produce consistent estimation rules $\{p_n(y \mid x)\}_{n=1}^{\infty}$ of the posterior $P(Y = y \mid X = x)$ which we denote as $p(y \mid x)$. It is important to note that consistency proof is independent of the splitting mechanism at each leaf node and so switching to the Gini impurity score and away from the mean squared error score of the GRF only has the potential to affect convergence rates.

To show that the posterior probability estimates are consistent, we need to show that $p_n(y \mid x) \xrightarrow{P} p(y \mid x)$ as $n \to \infty$. Let $y$ be the fixed arbitrary class label of interest. Given our data, we seek to estimate the class-specific posterior probability $p(y \mid x)$, equivalent to estimating the conditional mean $\theta(x) := \mathbb{E}[\mathbb{I}[Y = y] \mid X = x]$ for any $y \in \mathcal{Y}$. For conciseness, in the following proofs we will often drop the notational dependence of $\theta$ on $x$ and $y$ where convenient, letting it implicitly be a function of any fixed pair $(x, y) \in \mathcal{X} \times \mathcal{Y}$. To follow the notation of Athey et al. [29], we frame $\theta(x)$ as the solution to the estimation equation

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

where the score function $\psi_\theta(Y)$ is defined as

$$\psi_\theta(Y) := \mathbb{I}[Y = y] - \theta(x).$$

This estimand $\theta(x)$ can be estimated by solution $\hat{\theta}(x)$ to the empirical estimation equation

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

It follows that

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

per the expansion

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

$$= \sum_{i=1}^{n} \alpha_i(x)\mathbb{I}[Y_i = y] - \hat{\theta}(x) = 0.$$

These weights we learn from a learned random forest. Let a forest be composed of $B$ trees. In a single tree $b$, let $l_b(x)$ denote the set of training examples at the leaf node for which $X$ is placed. Define the weights $\alpha_{ib}(x)$ for that tree as

$$\alpha_{ib}(x) := \frac{1}{|l_b(x)|}\mathbb{I}[x_i \in l_b(x)],$$

the normalized indicator of whether or not $x$ and $x_i$ exist in the same leaf. Thus the forest weights $\alpha_i(x) = \frac{1}{B}\sum_{b=1}^{B} \alpha_{ib}(x)$ are simply the normalized weights across all trees.

By Theorem 3 of Athey et al. [29], a random forest built according to Specification 1 and solving an estimation problem satisfying Assumptions 1A-6A yields a consistent estimator. We enumerate these assumption, defined above, and verify that they hold for posterior probability estimate $\hat{\theta}(x)$.

1A. This remains a true assumption on the distribution and so is restated as Assumption 2.

2A. This is true, as evident in the derivatives

$$\frac{\partial}{\partial \theta}M_\theta(x) = -1 \quad \text{and} \quad \frac{\partial^2}{\partial^2 \theta}M_\theta(x) = 0.$$

3A. This is evident in the worse-case variogram for two solutions $\theta(x)$ and $\theta'(x)$

$$\gamma(\theta(x), \theta'(x))$$
$$:= \sup_{x \in \mathcal{X}} \{Var(\psi_\theta(Y) - \psi_{\theta'}(Y) \mid X = x)\}$$
$$= \sup_{x \in \mathcal{X}} \{Var(\theta'(x) - \theta(x)|X = x)\} = 0$$

which is trivially Lipschitz-continuous.

4A. Clearly $\psi_\theta(Y)$ is linear in $\theta(x)$ and so is a Lipschitz-continuous function in $\theta(x)$. The other term is 0 in this case.

5A. As shown previously shown, the estimation equation is solved to equal 0.

6A. This holds true by construction of the convex function $\Psi_\theta(Y) := \frac{1}{2}(\mathbb{I}[Y = y \mid X = x] - \theta(x))^2$ such that $\psi_\theta(Y) = -\frac{d}{d\theta}\Psi_\theta(Y)$, and the strongly convex function $\mathbb{M}_\theta(x) := \frac{1}{2}(P(Y = y \mid x) - \theta(x))^2$ where $M_\theta(x) = -\frac{d}{d\theta}\mathbb{M}_\theta(x)$.

This verifies Assumptions 2A-6A from Theorem 1 for the finite-sample estimate $\hat{p}_n(y|x) := \hat{\theta}(x)$ of $\theta(x) := p(y|x)$. Corollary 2 follows, adding Assumption 1A to the required Specifications 1-2 from Theorem 1.

*Proof of Theorem 3: A consistent classification rule..* Corollary 2 established consistency for each posterior probability estimate $p_n(y \mid x)$. We now proceed to show consistency for the classification rule $g_n(x) = \text{argmax}_y\, p_n(y \mid x)$. As before, define $p(y \mid x) := P(Y = y|X = x)$

LEMMA 2. *Let $x \in \mathcal{X}$ with true, but unknown, unique maximum $y^* := \operatorname{argmax}_y p(y \mid x)$, and define the finite sample estimate $\hat{y} := \operatorname{argmax}_y p_n(y \mid x)$. If $p_n(y \mid x)$ is a consistent estimator for $p(y \mid x)$, then*

$$P[\hat{y} \neq y^* \mid x] \to 0 \quad as \quad n \to \infty$$

*Proof.* We omit the conditional for notational brevity by substituting $p(y) := p(y \mid x)$ and $p_n(y) := p_n(y \mid x)$. Then it follows that

$$P[\hat{y} \neq y^* \mid x] = P[\max_y p_n(y) > p_n(y^*)]$$

$$= P\left[\bigcup_{y \neq y^*} p_n(y) > p_n(y^*)\right]$$

$$\leq \sum_{y \neq y^*} P[p_n(y) > p_n(y^*)]$$

$$= \sum_{y \neq y^*} P[p_n(y) - p_n(y^*) > 0]$$

$$= \sum_{y \neq y^*} P[(p_n(y) - p_n(y^*)) -$$

$$(p(y) - p(y^*)) > p(y^*) - p(y)]$$

Let $\varepsilon_y := p(y^*) - p(y)$ and note that $\varepsilon_y > 0$ for all $y \in \mathcal{Y} \setminus \{y^*\}$ since $y^*$ is a unique maximum. Observe that

$$\sum_{y \neq y^*} P[(p_n(y) - p_n(y^*)) - (p(y) - p(y^*)) > \varepsilon_y]$$

$$\leq \sum_{y \neq y^*} P\left[\left|(p_n(y) - p_n(y^*)) - (p(y) - p(y^*))\right| > \varepsilon_y\right]$$

By the consistency of the individual posteriors, the difference of two is consistent and so since $\mathcal{Y}$ is a finite set,

$$P[\hat{y} \neq y^* \mid x]$$

$$\leq \sum_{y \neq y^*} P[|(p_n(y) - p_n(y^*)) - (p(y) - p(y^*))| > \varepsilon_y]$$

$$\to 0 \quad as \quad n \to \infty \qquad \square$$

With Lemma 2, the proof of Theorem 3 follows. Denote the finite samples classification rule $\hat{y} := \operatorname{argmax}_y p_n(y \mid x)$ as before and let $y^* := \operatorname{argmax}_y p(y \mid x)$ be a unique maximum. If $y^*$ were not unique, we would instead consider the aggregate of all such maximum classes as a pseudo class, apply the following analyses, and be confident in both $L_n$ and $L^*$ up to a factor equal to the reciprocal of the number aggregated classes due to a chance guess between them.

Otherwise, for any $\varepsilon > 0$, by the law of total probabilities,

$$P\left[|L_n - L^*| > \varepsilon\right]$$
$$= P\left[|p_n(\hat{y} \mid x) - p(y^* \mid x)| > \varepsilon\right]$$
$$= P\left[|p_n(\hat{y} \mid x) - p(y^* \mid x)| > \varepsilon \mid \hat{y} = y^*\right] \times P\left[\hat{y} = y^*\right]$$
$$+ P\left[|p_n(\hat{y} \mid x) - p(y^* \mid x)| > \varepsilon \mid \hat{y} \neq y^*\right] \times P\left[\hat{y} \neq y^*\right].$$

In the case that $\hat{y} = y^*$, by Corollary 2 we have convergence of the posteriors and so

$$P\left[|p_n(\hat{y} \mid x) - p(y^* \mid x)| > \varepsilon \mid \hat{y} = y^*\right] \to 0 \quad \text{as } n \to \infty.$$

In the case that $\hat{y} \neq y^*$, by Lemma 2 we have that

$$P\left[\hat{y} \neq y^*\right] \to 0 \quad \text{as } n \to \infty.$$

Since the probabilities are bounded above by one, it follows that as $n \to \infty$,

$$P\left[|p_n(\hat{y} \mid x) - p(y^* \mid x)| > \varepsilon \mid \hat{y} = y^*\right] \times P\left[\hat{y} = y^*\right] \to 0$$

and

$$P\left[|p_n(\hat{y} \mid x) - p(y^* \mid x)| > \varepsilon \mid \hat{y} \neq y^*\right] \times P\left[\hat{y} \neq y^*\right] \to 0$$

and thus

$$P\left[|L_n - L^*| > \varepsilon\right] = P\left[|p_n(\hat{y} \mid x) - p(y^* \mid x)| > \varepsilon\right] \to 0$$

**Appendix B. Experiment Extras.**

**B.1. Mathematical description of sampling manifolds in simulation examples.** In 3.2.1

**B.2. Three simulated manifolds: time complexity.** Each simulated experiment was run on CPUs and allocated 52 cores for parallel processing. The resulting train and test times as a function of the number of training samples are plotted in Figure S1. MORF has train and test times on par with those of SPORF and so is not particularly more computationally intensive to run. The ConvNet, however, took noticeably longer to run across simulations for the majority of sample sizes.



Fig. S1: Algorithm train times (above) and test times (below) across increasing sample sizes. MORF runtime is not particularly costly and well below ConvNet runtime in most examples.

**B.3. Intracranial EEG Experiments - Ethics.** For the motor control, details of the experiment are in Kerr et al. [37]. If the patient expressed interest in participating, the research staff would verbally review the written, IRB approved consent form. If agreed upon, the patient would sign the written consent and be enrolled in the study. A copy of the written consent would also be given to patient to keep.

Experimental protocols were approved by the Cleveland Clinic Institutional Review Board.

For the epilepsy intracranial EEG data, details on the dataset can be found at [42]. All data were acquired with approval from the local institutional review board (IRB) at each clinical institution: UMMC by the IRB of the University of Maryland School of Medicine; UMH by the University of Miami Human Subject Research Office—Medical Sciences IRB; NIH by the National Institutes of Health IRB; JHH by Johns Hopkins IRB; and CClinic by the Cleveland Clinic IRB. Informed consent was given at each clinical center. The acquisition of data for research purposes was completed with no impact on the clinical objectives of the patient stay. Digitized data were stored in an IRB-approved database compliant with Health Insurance Portability and Accountability Act regulations.

### B.4. 1D Locality: Predicting Movement Direction With Intracranial EEG.

MORF's performance was next evaluated on stereotactic electroencephalogram (sEEG) data recorded in epilepsy patients undergoing a motor control task presented in Kerr et al. [37], Breault et al. [48]. The classification task presented here is to predict movement direction (up, down, left, or right) based on the sEEG data alone, rather then performing explicit feature engineering, such as computing power in frequency bands. We compare MORF to other classification algorithms. The interesting aspect of this data is that there are no motor regions recorded. Thus, our hypothesis is that only a subset of the recording electrodes over time are important in decoding movement directionality. This is analogous to Experiment D mentioned in Section 5.2. Each subject performed the task for several trials, each consisting of a movement instruction followed by a movement generated by the subject. Using only the sEEG data, we sought to decode the movement directionality. For full details on the dataset and clinical problem, we refer the readers to Kerr et al. [37], Breault et al. [48].

Here we perform 5-fold cross validation for each subject including all the sEEG recording electrodes time-locked to a movement onset marking. Each fold is *a priori* generated per subject, where there is a set of testing trials left out. The overall task is very challenging because there are no motor brain regions being recorded. Nonetheless, we expect that other brain regions are involved in the motor control process. MORF is able to achieve a superior performance measured by AUC relative to the other classifiers, as seen in Figure S2. Notably, MORF and SPORF perform the best in this setting with a limited set of training samples, whereas the ConvNet performs slightly better then chance on the test set, overfitting to training set. Across the set of all folds of all subjects, MORF was never worse than another classifier in terms of the median pairwise difference in Cohen's kappa while being significantly (p-value $\leq 0.05$) better than the MLP and ConvNet per a Wilcoxon paired sign test on those same pairwise differences.

### B.5. Three simulated manifolds: model-misspecification.
For each simulation problem, we optimized the optimal parameters for MORF on a large dataset of 2000 samples using grid search over a range of valid parameter values for the manifold structure. Once we arrived at an optimal parametrization, we then proceeded to modify the parametrizations for each of the experiments to change the dimensions of the possible patches MORF could sample. For full details, see the online repository, where the experiment was done (https://github.com/adam2392/morf-demo). We then performed 10-fold stratified cross-validation for each parameter setting, computing accuracy on the held-out test set. We randomly produced 500 samples from each simulated
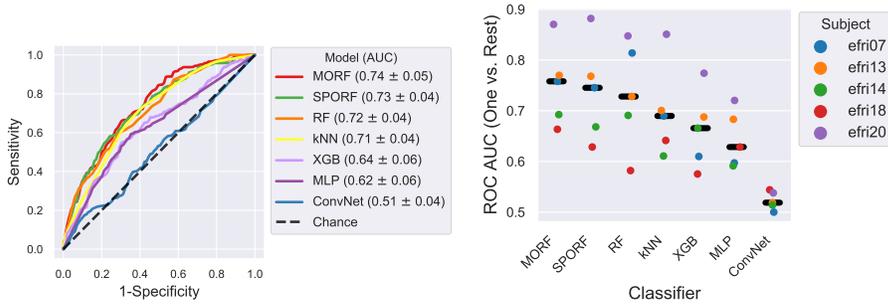
Fig. S2: MORF performance on decoding movement direction from the raw sEEG data in non-motor brain regions. Subjects are undergoing a motor-control task. The naming of subjects is simply derived from their clinical monitoring session and does not reflect any specific numbering scheme. **(Top)** Shows ROC curve of moving down in the motor task decoded with all classifiers on the same set of data and their AUC scores. **(Bottom)** Shows a summary AUC stripplot where each dot represents the held-out trial median AUC score for a certain subject over 5-fold CV, and the median of the overall AUC for each classifier is shown (solid black line). In almost all subjects, MORF gains in AUC compared to the other classifiers with fixed hyperparameters and fixed trials in each of the 5 folds.
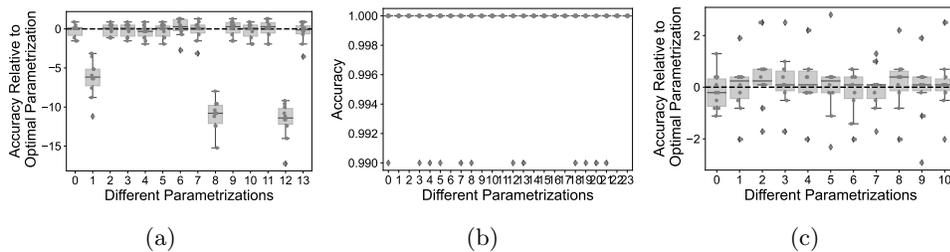


Fig. S3: MORF performance when model is mis-specified. (a) Circle segments, (b) Horizontal bars, and (c) a 1D time-series with an impulse. The dashed line in (a) and (c) indicate a consistent accuracy score relative to the optimal parametrizations. Each simulation (b) is shown with just accuracy because the relatively low variance in the optimal parameter classifier made the resulting normalized plot uninterpretable. The low performance of a few parameterizations in (a) were when the maximum possible patch size was set to 6. This likely restricts MORF from learning the Circle Segment structure as fast as when MORF can sample larger patches.

dataset, and kept all other MORF parameters the same. We used a total of 500 trees for each problem. Each non-optimal parametrization was normalized with respect to the accuracy scores of the optimal parametrization (subtracting the mean and dividing by the standard deviation of the optimal scores). The normalized accuracy scores of each non-optimal parameter setting are then shown in Supplementary Figure S3.

28

**Appendix C. Pseudocode.**

---

**Algorithm 1** Learning a Manifold Oblique decision tree, modified from Tomita et al. [9].

---

**Input:** (1) $\mathcal{D}_n$: training data (2) $d$: dimensionality of the projected space, (3) $f_{\mathbf{A}}$: distribution of the atoms, (4) $\Theta$: set of split eligibility criteria
**Output:** A MORF decision tree $T$
1: **function** $T = \text{GROWTREE}(\mathbf{X}, \mathbf{y}, f_{\mathbf{A}}, \Theta)$
2:     $c = 1$                            $\triangleright$ $c$ is the current node index
3:     $M = 1$                     $\triangleright$ $M$ is the number of nodes currently existing
4:     $S^{(c)} = \text{bootstrap}(\{1, ..., n\})$ $\triangleright$ $S^{(c)}$ is the indices of the observations at node $c$
5:     **while** $c < M + 1$ **do**               $\triangleright$ visit each of the existing nodes
6:         $(\mathbf{X}', \mathbf{y}') = (\mathbf{x}_i, y_i)_{i \in S^{(c)}}$               $\triangleright$ data at the current node
7:         **for** $k = 1, \ldots, K$ **do** $n_k^{(c)} = \sum_{i \in S^{(c)}} I[y_i = k]$ **end for**   $\triangleright$ class counts (for classification)
8:         **if** $\Theta$ satisfied **then**                $\triangleright$ do we split this node?
9:             $\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_d] \sim f_{\mathbf{A}}$        $\triangleright$ sample random $p \times d$ matrix of atoms
10:            $\widetilde{\mathbf{X}} = \mathbf{A}^T \mathbf{X}' = (\widetilde{\mathbf{x}}_i)_{i \in S^{(c)}}$    $\triangleright$ random projection into new feature space
11:            $(j^*, t^*) = \text{findbestsplit}(\widetilde{\mathbf{X}}, \mathbf{y}')$              $\triangleright$ Algorithm 2
12:            $S^{(M+1)} = \{i : \mathbf{a}_{j^*} \cdot \widetilde{\mathbf{x}}_i \leq t^* \quad \forall i \in S^{(c)}\}$      $\triangleright$ assign to left child node
13:            $S^{(M+2)} = \{i : \mathbf{a}_{j^*} \cdot \widetilde{\mathbf{x}}_i > t^* \quad \forall i \in S^{(c)}\}$    $\triangleright$ assign to right child node
14:            $\mathbf{a}^{*(c)} = \mathbf{a}_{j^*}$             $\triangleright$ store best projection for current node
15:            $\tau^{*(c)} = t^*$            $\triangleright$ store best split threshold for current node
16:            $\kappa^{(c)} = \{M + 1, M + 2\}$        $\triangleright$ node indices of children of current node
17:            $M = M + 2$           $\triangleright$ update the number of nodes that exist
18:         **else**
19:            $(\mathbf{a}^{*(c)}, \tau^{*(c)}, \kappa^{*(c)}) = \text{NULL}$
20:         **end if**
21:         $c = c + 1$                    $\triangleright$ move to next node
22:     **end while**
23:     **return** $(S^{(1)}, \{\mathbf{a}^{*(c)}, \tau^{*(c)}, \kappa^{(c)}, \{n_k^{(c)}\}_{k \in \mathcal{Y}}\}_{c=1}^{m-1})$
24: **end function**

---

29

**Algorithm 2** As in Tomita et al. [9]. Finding the best node split. This function is called by growtree (Alg 1) at every split node. For each of the $p$ dimensions in $\mathbf{X} \in \mathbb{R}^{p \times n}$, a binary split is assessed at each location between adjacent observations. The dimension $j^*$ and split value $\tau^*$ in $j^*$ that best split the data are selected. The notion of "best" means maximizing some choice in scoring function. In classification, the scoring function is typically the reduction in Gini impurity or entropy. The increment function called within this function updates the counts in the left and right partitions as the split is incrementally moved to the right.

---

**Input:** (1) $(\mathbf{X}, \mathbf{y}) \in \mathbb{R}^{p \times n} \times \mathcal{Y}^n$, where $\mathcal{Y} = \{1, \ldots, K\}$
**Output:** (1) dimension $j^*$, (2) split value $\tau^*$

1: **function** $(j^*, \tau^*) = \text{FINDBESTSPLIT}(\mathbf{X}, \mathbf{y})$
2:     **for** $j = 1, \ldots, p$ **do**
3:         Let $\mathbf{x}^{(j)} = (x_1^{(j)}, \ldots, x_n^{(j)})$ be the $j$th row of $\mathbf{X}$.
4:         $\{m_i^j\}_{i \in [n]} = \text{sort}(\mathbf{x}^{(j)})$   $\triangleright$ $m_i^j$ is the index of the $i^{th}$ smallest value in $\mathbf{x}^{(j)}$
5:         $t = 0$                    $\triangleright$ initialize split to the left of all observations
6:         $n' = 0$                $\triangleright$ number of observations left of the current split
7:         $n'' = n$              $\triangleright$ number of observations right of the current split
8:         **if** (task is classification) **then**
9:             **for** $k = 1, \ldots, K$ **do**
10:                 $n_k = \sum_{i=1}^n I[y_i = k]$     $\triangleright$ total number of observations in class $k$
11:                 $n_k' = 0$  $\triangleright$ number of observations in class $k$ left of the current split
12:                 $n_k'' = n_k$     $\triangleright$ number of observations in class $k$ right of the current split
13:             **end for**
14:         **end if**
15:         **for** $t = 1, \ldots, n - 1$ **do** $\triangleright$ assess split location, moving right one at a time
16:             $(\{(n_k', n_k'')\}, n', n'', y_{m_t^j}) = \text{increment}(\{(n_k', n_k'')\}, n', n'', y_{m_t^j})$
17:             $Q^{(j,t)} = \text{score}(\{(n_k', n_k'')\}, n', n'')$         $\triangleright$ measure of split quality
18:         **end for**
19:     **end for**
20:     $(j^*, t^*) = \underset{j,t}{\text{argmax}}\, Q^{(j,t)}$
21:     **for** $i = 0, 1$ **do** $c_i = m_{t^*+i}^{j^*}$ **end for**
22:     $\tau^* = \frac{1}{2}(x_{c_0}^{(j^*)} + x_{c_1}^{(j^*)})$    $\triangleright$ compute the actual split location from the index $j^*$
23:     **return** $(j^*, \tau^*)$
24: **end function**

**Appendix D. Hyperparameters.**

Table S1: ConvNet hyperparameters for each experiment.

| Experiment | Classifier | Architecture Sequence |
|---|---|---|
| Circle | ConvNet | Conv1d(32, window=6, stride=1) |
| | | MaxPool1d(window=2, stride=2) |
| | | Conv1d(64, window=10, stride=1) |
| | | MaxPool1d(window=2, stride=2) |
| | | Dropout(p=0.5), Linear(500, 2) |
| H/V Bars | ConvNet | Conv2d(32, window=5, stride=1) |
| | | MaxPool1d(window=2, stride=2) |
| | | Conv2d(64, window=5, stride=1) |
| | | MaxPool1d(window=2, stride=2) |
| | | Dropout(p=0.5), Linear(200, 2) |
| Impulse | ConvNet | Conv1d(32, window=10, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(64, window=5, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Dropout(p=0.5), Linear(200, 2) |
| Experiments D, E | ConvNet | Conv2d(32, window=2, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(32, window=5, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(64, window=5, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Linear(64), Linear(n_classes) |
| MNIST | ConvNet | Conv2d(32, window=5, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(64, window=5, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Dropout(p=0.5), Linear(200, 10) |
| Surgical Outcome | ConvNet | Conv2d(32, window=2, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(32, window=5, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(64, window=5, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Linear(64), Linear(n_classes) |
| Predicting Movement | ConvNet | Conv2d(32, window=3, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(64, window=3, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Conv2d(64, window=3, stride=1) |
| | | MaxPool2d(window=2, stride=2) |
| | | Linear(64), Linear(n_classes) |

Table S2: *scikit-learn*, SPORF, and MORF hyperparameters for each experiment.

| Experiment | Classifier | Hyperparameters |
|---|---|---|
| All | Lin. SVM | C=1, penalty="l2";kernel="linear";loss="squared_hinge" |
| All | Log. Reg | C=1, penalty="l2" |
| All | MLP | activation="relu"; alpha=0.0001; hidden_layer_sizes=(100,); solver="adam" |
| All | RF | n_trees=500, max_features='sqrt' |
| All | XGB | n_boosting_rounds = 10; learning_rate=0.3; max_depth=6; subsample=1; tree_method='auto' (all default) |
| All | SPORF | n_trees=500, max_features='sqrt' |
| All | SVM | C=1; gamma=1/(n_features*Var(X)); kernel="rbf" |
| All | kNN | n_neighbors=5; p=2 |
| Circle | MORF | n_trees=500, max_features=0.5; patch_height_max=1; patch_height_min=1; patch_width_max=12; patch_width_min=3 |
| H/V Bars | MORF | n_trees=500, max_features='sqrt'; patch_height_max=2; patch_height_min=2; patch_width_max=9; patch_width_min=2 |
| Impulse | MORF | n_trees=500, max_features=0.3; patch_height_max=1; patch_height_min=1; patch_width_max=12; patch_width_min=2 |
| Experiment D | MORF | n_trees=500, max_features='sqrt'; patch_height_max=2; patch_height_min=2; patch_width_max=10; patch_width_min=5; |
| Experiment E | MORF | n_trees=500, max_features='sqrt'; patch_height_max=2; patch_height_min=1; patch_width_max=20; patch_width_min=5; |
| MNIST | MORF | n_trees=500, max_features='sqrt' patch_height_max=2; patch_height_min=2; patch_width_max=5; patch_width_min=2; |
| Surgical Outcome | MORF | n_trees=500, max_features='sqrt'; patch_height_max=sqrt(height); patch_height_min=1; patch_width_max=sqrt(width); patch_width_min=1; |