
Randomer Forests

Abstract

1. Introduction

Data science is becoming increasingly important as our ability to collect and process data continues to increase. Supervised learning—the act of using predictors to make a prediction about some target data—is of special interest to many applications, ranging from science to government to industry. Classification, a special case of supervised learning, in which the target data is categorical, is one of the most fundamental learning problems, and has been the subject of much study. A simple pubmed search for the term “classifier” reveals nearly 9,000 publications, and a similar arXiv search reports that the number of hits is >1000. One of the most popular and best performing classifiers is random forests (RFs) (Breiman, 2001). Several recent benchmark papers assess the performance of many different classification algorithms on many different datasets (Fernandez-Delgado et al., 2014; Caruana et al., 2008), and both concluded the same thing: RFs are the best classifier.

RF typically refers to Breiman’s Forest-IC, which uses axis-parallel (Heath et al., 1993), or orthogonal trees (Menze et al., 2011). That is, the feature space is recursively split along directions parallel to the axes of the feature space. Thus, in cases in which the class marginal distributions seem inseparable, RF may be suboptimal. To address this, Breiman also proposed and characterized Forest-RC, which used linear combinations of coordinates rather than individual coordinates, to split along. Others have studied several different variants of “oblique” decision forests, including efforts to learn good projections (Heath et al., 1993; Tan & Dowe, 2005), using principal components analysis to find the directions of maximal variance (Ho, 1998; Rodriguez et al., 2006), or directly learning good discriminant directions (Menze et al., 2011). While all of these recent approaches deal with rotational invariance, they fail to address two important issues. First, in real data, the number of *irrelevant* features tends to be large, especially for high dimensional data, giving rise to optimal

decision boundaries that are sparse. In such cases, we say that the signal is *compressive*. To wit, all of the proposed oblique decision forests, with the exception of Forest-RC if tuned appropriately, do not impose a constraint on the sparsity of recursive partitions of the feature space. Therefore, such methods may be suboptimal in cases in which the signal is compressive. Second, as our ability to collect data and process continues to increase, we are encountering increasingly massive datasets. Therefore, time- and space-efficient methods are becoming more and more important. With the exception of Forest-RC, all oblique methods require expensive computation and/or storage. There is therefore a gap that calls for the construction of a scalable oblique decision forest classifier that can perform well in the presence of an overwhelming number of irrelevant features.

To bridge this gap, we first state a generalized forest building scheme, random projection forests, which includes all of the above algorithms as special cases. This enables us to formally state our objective, and provides a lens that enables us to propose a few novel special cases, which we refer to as “*Randomer Forests*” (or RerFs for short), for reasons that will become clear in the sequel. We both theoretically and empirically demonstrate that our methods have the same time and space complexity as random forests, and show on simulated datasets that it performs well in the face of many irrelevant features and when axis-aligned splits are suboptimal. Additionally, we demonstrate that our method empirically outperforms both RF and random rotation RF (RotRF), which is another oblique method (Blaser & Fryzlewicz, 2015), on a suite of benchmark datasets. Lastly, we propose two augmentations to RerF for making the method more robust to affine transformations and outliers. To conclude, we propose our method as a competitive alternative to RF and other oblique decision forest methods. Open source code will be made available.

2. Random Projection Forests

Let $\mathcal{D}^n = \{(X_i, y_i) : i \in [n]\}$ be a given dataset, where $X_i \in \mathbb{R}^p$ and $y_i \in \mathcal{Y} = \{\tilde{y}_1, \dots, \tilde{y}_C\}$. A classification forest, $\bar{g}(X; \mathcal{D}^n)$ is an ensemble of L decision trees, each tree $g^l(X; \mathcal{D}^l)$ is trained on a (sub)set of the data, $\mathcal{D}^l \subset \mathcal{D}^n$. Random Projection forests are a special case of classification forests that subsume all of the strategies mentioned above (see Pseudocode 1). The key idea of all

of them is that at each node of the tree, we have a set of predictor data points, $\bar{X} = \{X_s\}_{s \in \mathcal{S}_{ij}^l} \in \mathbb{R}^{p \times S_{ij}^l}$, where $S_{ij}^l = |\mathcal{S}_{ij}^l|$ is the cardinality of the set of predictor data points at the $(ij)^{th}$ node of the l^{th} tree. We sample a matrix $A \sim f_A(\mathcal{D}^n)$, where $A \in \mathbb{R}^{p \times d}$, possibly in a data dependent fashion, which we use to project the predictor matrix \bar{X} onto a lower dimensional subspace, yielding $\tilde{X} = A^T \bar{X} \in \mathbb{R}^{d \times S_{ij}^l}$, where $d \leq p$ is the dimensionality of the subspace. To wit, Breiman's original Forest-IC algorithm can be characterized as a special case of random projection forests. In particular, in Forest-IC one constructs A such that for each of the d columns, we sample a coordinate (without replacement), and put a 1 in that coordinate, and zeros elsewhere. Breiman's Forest-RC constructs A by sampling a fixed number of coordinates (without replacement) for each of the d columns, and puts a value uniformly sampled from $[-1, 1]$ in each of those coordinates. Similarly, Ho's rotation forests construct A from the top d principal components of the data \bar{X} at a given node. Thus, the key difference in all these approaches is the choice of f_A .

Algorithm 1 Psuedocode for Random Projection Forests, which generalizes a wide range of previously proposed decision forests.

```

Input: data:  $\mathcal{D}^n = (X_i, y_i) \in (\mathbb{R}^p \times \mathcal{Y})$  for  $i \in [n]$ , tree
rules (stopping criteria, rules for sampling data points
per tree, etc.), distributions on  $d \times p$  matrices:  $A \sim
f_A(\mathcal{D}^n)$ , preprocessing rules
Output: decision trees, predictions, out of bag errors,
etc.
Preprocess according to rule
for each tree do
    Subsample data to obtain  $(\bar{X}, \bar{y})$ , the set of data points
    to be used in this tree
    for each leaf node in tree do
        Let  $\tilde{X} = A^T \bar{X} \in \mathbb{R}^{d \times s}$ , where  $A \sim f_A(\mathcal{D}^n)$ 
        Find the “best” split coordinate  $k^*$  in  $\tilde{X}$  and “best”
        split value  $t^*(k^*)$  for this coordinate
        Split  $X$  according to whether  $X(k) > t^*(k^*)$ 
        Assign each child node as a leaf or terminal node
        according to stopping criteria
    end for
end for
Prune trees according to rule

```

The goal of this work is to find a random projection forest that possesses many of the desirable properties of RF, yet is able to find decision boundaries that are less biased by geometrical constraints (i.e. not constrained to be axis-aligned), in part by changing the distribution f_A . The result we call randomer forests (or RerFs for short).

3. Randomer Forests

Our choice in f_A is based on the following three ideas:

1. While RF empirically performs well in many settings, it is quite restrictive in that candidate splits evaluated at each node are constrained to be axis-aligned. It is sometimes the case that axis-aligned splits are suboptimal, and oblique splits may be desired.
2. It is often the case that the number of *relevant* features is substantially smaller than the number of total features. We say that the signal is *compressive* in such cases. RF often does well when the signal is compressive, which may be attributed, in part, to the extreme sparsity constraint on the columns of A . All of the oblique decision forest algorithms to date do not impose a sparsity constraint on the candidate split directions, which may cause such algorithms to underperform when the signal is compressive.
3. Except for Breiman’s Forest-RC, existing oblique decision forest algorithms involve expensive computations to identify and select splits, rendering them less space and time efficient than Forest-IC or Forest-RC. An oblique decision forest having a space and time complexity comparable to Forest-IC or Forest-RC is desirable.

To this end, we employ very **sparse random projections** (Li et al., 2006). Rather than sampling d non-zero elements of A and enforcing that each column gets a single non-zero number (without replacement) as RF does, we relax these constraints and select d non-zero numbers from $\{-1, +1\}$ with equal probabilities and distribute uniformly at random in A .

It is apparent that RerF bares a resemblance to Forest-RC, yet there is one key difference. In Forest-RC, the number of nonzeros in each column of A is fixed to be the same across A and for every split node, and its optimal value has to be found through parameter selection. Our construction circumvents selection of this parameter by randomizing the number of nonzeros in each column of A . Furthermore, our algorithm allows A to have columns of varying sparsity, which may promote more diversity in the ensemble. To our knowledge, this property does not exist in any of the proposed random projection forest algorithms.

4. Experimental Results

4.1. Simulations Involving Compressive Signals

Many classification problems arise in which the signal is compressive and the optimal split directions are not axis-aligned. We constructed two synthetic datasets with both

110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

220 of these properties to compare classification performance
 221 and training time of RF, RerF, and RotRF:

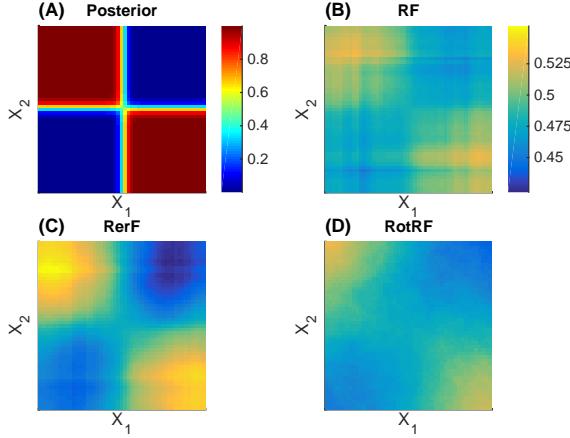
222 **Sparse parity** is a variation of the parity problem. The
 223 parity problem is a multivariate generalization of the XOR
 224 problem and is one of the hardest constructed binary clas-
 225 sification problems. In parity, a given sample has a mean
 226 whose elements are Bernoulli samples with probability 1/2,
 227 and then Gaussian noise is independently added to each di-
 228 mension with the same variance. A sample's class label
 229 is equal to the parity of its mean. Sparse parity is an adap-
 230 tion of the basic parity problem in which the sample's
 231 class label is equal to the parity of only the first p^* elements
 232 of the mean, rendering the remaining $p - p^*$ dimensions as
 233 noise.

234 **Trunk** is a well-known binary classification in which each
 235 class is distributed as a p -dimensional multivariate Gaus-
 236 sian with identity covariance matrices (Trunk, 1979). The
 237 means of the two classes are $\mu_1 = (1, \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{3}}, \dots, \frac{1}{\sqrt{p}})$ and
 238 $\mu_2 = -\mu_1$. It follows from this that the signal-to-noise ratio
 239 of the i th dimension asymptotically decreases to zero
 240 with increasing i .

241 RotRF is an oblique decision forest method that has been
 242 recently proposed in (Blaser & Fryzlewicz, 2015). RotRF
 243 works exactly like RF, except that the feature space is uni-
 244 formly randomly rotated prior to inducing each decision
 245 tree. This induces oblique splits in an unstructured way,
 246 which the authors argue is favorable to structured methods
 247 such as in Rodriguez's rotation forest, because such struc-
 248 tured methods reduce the diversity of the trees. Uniformly
 249 random rotations of the feature space imply that in gen-
 250 eral, splits will not be sparse. Therefore, we conjecture
 251 that RotRF will perform increasingly poorly as the ratio of
 252 the number of irrelevant features to the number of relevant
 253 features becomes larger, while RF and RerF will be rela-
 254 tively more robust to the increasing presence of irrelevant
 255 features.

256 Figure 1 depicts the true class posteriors in the first two
 257 dimensions of the sparse parity simulation in panel A and
 258 estimates of the posteriors for RF, RerF, and RotRF in pan-
 259 el B, C, and D, respectively. For this simulation, $p = 20$,
 260 $p^* = 3$, and $n = 1000$. The number of trees used for all
 261 three algorithms was 1000. Various values of d , the number
 262 of columns in the random projection matrix A , were tried
 263 and the best for each algorithm was selected. RerF gives
 264 estimates of the posteriors closest to the true posteriors.

265 The left panels of Figure 2 show two-dimensional scatter
 266 plots from each of the two example simulations (using
 267 the first two coordinate dimensions). The middle panels
 268 show the misclassification rate relative to RF against the
 269 number of observed dimensions p , for RerF and RotRF.
 270 Relative misclassification rate was computed as the differ-



275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329

Figure 1. Class posteriors for the sparse parity problem in the first two dimensions (refer to section 4.1 for details). (A): True posteriors. (B) - (D): Posterior estimates for RF, RerF, and RotRF, respectively. Comparing panels B-D with A, it is evident that RerF produces better estimates of the posteriors than does RF or RotRF

ence between the misclassification rate of either RerF or RotRF and that of RF. The misclassification rate of RF relative to itself is shown for reference. The right panels show training time against the number of observed dimensions p for all four classifiers. For all comparisons, we used the same number of trees for each method to enable a fair comparison. The number of trees used for each method in the sparse parity and Trunk simulations were 500 and 1000, respectively. In all methods, trees were pruned and the minimum number of data points at a node in order to be considered for splitting was 10. Gini impurity was used as the split criteria. The only parameter tuned was d , the number of candidate split directions evaluated at each split node. When $p \leq 5$, each classifier was trained for all $d \in [p]$. When $p > 5$, each classifier was trained for all $d \in \{1, p^{1/4}, p^{1/2}, p^{3/4}, p\}$. Misclassification rates for each classifier were selected as the lowest achieved from the different values of d tried. Training times for each classifier were computed by averaging the training times given by all values of d tried. For sparse parity, n was fixed at 1000 and classifiers were evaluated for $p \in \{2, 5, 10, 25, 50, 100\}$. The relevant number of features p^* was fixed at a value of 3. For Trunk, n was fixed at 100 and classifiers were evaluated for $p \in \{2, 10, 50, 100, 500\}$. Relative error and training times plotted are the average of 25 repeated trials, and error bars represent the standard error of the mean.

In panel B, RerF performs as well as or better than both RF and RotRF for all values of p . RotRF performs as well as

or better than RF except for when $p = 25$. As conjectured, RotRF performs better than RF when p is small because oblique splits provide an advantage over axis-aligned splits in the sparse parity problem. As p increases and the ratio $p * /p$ decreases, RerF begins to outperform RotRF. Ultimately, when this ratio is small enough, RotRF performs even worse than RF. RerF's ability to perform relatively well can be attributed to the sparsity of oblique splits. In panel E, RerF outperforms RF for all values of p . This is because linear combinations of a few features can yield a higher signal-to-noise ratio than any single feature. RotRF outperforms both RF and RerF up to $p = 100$. RotRF is able to perform better than RerF in these cases because a larger number of features can be linearly combined to yield an even higher signal-to-noise ratio. When $p = 500$, classification performance of RotRF significantly degrades and becomes worse than both RF and RerF. This can be explained by the fact that when p is large enough, RotRF often samples linear combinations of many features each having a low signal-to-noise ratio. Such projections will yield a lower signal-to-noise ratio than any single feature. In panels C and F, training times are comparable when p is small. As panel F indicates, training time of RotRF is significantly longer than those of RF and RerF. This is due to the fact that generating random rotation matrices requires QR decompositions, which have a time complexity of $O(p^3)$.

4.2. Theoretical Space and Time Complexity

For a random forest, assume there are L trees. If there are n data points per tree, and each tree grows until terminal nodes have only $\mathcal{O}(1)$ data points with d coordinates in them, there are $\mathcal{O}(n)$ nodes. Then the complexity of constructing the random forest, disregarding cross-validation or other randomization techniques for preventing overfitting, is $\mathcal{O}(Ln^2d \log n)$. In practice the trees are shallower and stop much earlier than when nodes have $\mathcal{O}(1)$ points, so “in practice” the complexity often appears to be $\mathcal{O}(Lnd \log n)$.

Randomer forest has a very similar time and space complexity, unlike many of the other oblique random forest variants. Specifically, assume that RerF also has L trees, and n data points per tree, and no pruning, so $\mathcal{O}(n)$ nodes. Like RF, RerF requires $\mathcal{O}(d)$ time to sample d non-zero numbers, and $\mathcal{O}(dn_k)$ time to obtain the new matrix, \tilde{X} , because it is a sparse matrix multiplication, in node k with $\mathcal{O}(n_k)$ points. RerF also takes another $\mathcal{O}(d/n \log(d/n)) = \mathcal{O}(1)$ time to find the best dimension. Thus, in total, in the case of well-balanced trees, RerF also requires only $\mathcal{O}(Ldn^2 \log n)$ time to train. To store the resulting RerF requires $\mathcal{O}(Ln \log n)$, because each node can be represented by the indices of the coordinates that received a non-zero element, and the expected number of such indices is $\mathcal{O}(1)$. The only additional space constraint is storing which in-

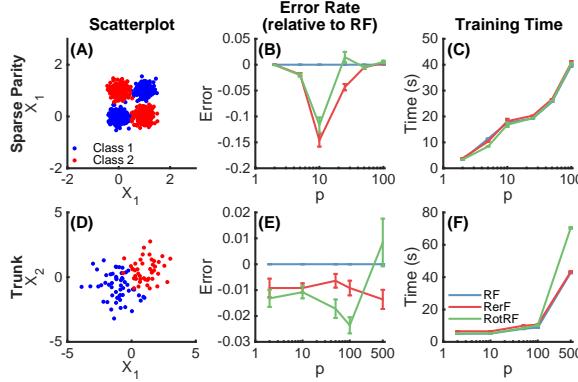


Figure 2. Sparse parity (A-C) and Trunk (D-F) simulations (see section 4.1 for details). (A) and (D): Scatterplots of sampled points in the first two dimensions. (B) and (E): Error rates of RerF and RotRF relative to RF across different values of p . (C) and (F): Same as (B) and (E) except absolute training time is plotted on the y-axis instead. The blue, red, and green lines correspond to RF, RerF, and RotRF, respectively. This color coding is used for all pertinent figures that follow. Both RerF and RotRF do better than RF when the number of irrelevant features is sufficiently small, due to their ability to generate oblique partitions. However, when the number of irrelevant features becomes large enough, performance of RotRF rapidly degrades. Training times show that RerF scales identically with RF, while RotRF scales poorly with large p .

dices are positive, and which are negative, which is merely another constant. The cost of initially sorting in RerF(r), and of computing the means of the classes in RerF(d) are negligible compared to the main cost above. Note that these numbers are in stark contrast to other oblique random forests. For example, rotation forests require running a singular value decomposition at each node.

4.3. Effects of Transformations and Outliers

We next want to compare the robustness of RF, RerF, and RotRF to various data transformations. To do so, we consider several different modifications to the simulation settings described in the previous section: rotation, scale, affine, and outliers. To rotate the data, we simply generate rotation matrices uniformly and apply them to the data. To scale, we applied a scaling factor sampled from a uniform distribution on the interval $[0,10]$ to each dimension. Affine transformations were performed by applying a combination of rotations and scalings as just described. Additionally, we examined the effects of introducing outliers. Outliers were introduced by sampling points from the distributions as previously described but instead using covariance matrices scaled up by a factor of four. Empirically, an

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

440

addition of 20 points from these outlier models to the original 100 points was found to produce a noticeable but not overwhelming effect on classifier performance.

Figure 3 shows the effect of these transformations and outliers on the sparse parity (panels A-E) and Trunk (panels F-J) problems. Comparing panels B and G with A and F show that RF and RerF are sensitive to rotations, while RotRF is insensitive to rotations. Panels C and H show that RF is insensitive to scale, while RerF is marginally affected and RotRF is significantly affected. Panels D and I show that all three algorithms are affected by affine transformations. Lastly, panels E and J demonstrate that all three algorithms are sensitive to outliers.

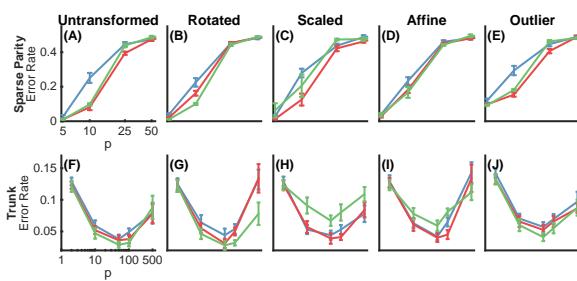


Figure 3. The effects of different transformations applied to the sparse parity (A-E) and Trunk (F-J) simulations on classification performance (see section 4.3 for details). Specifically, we consider rotations, scalings, affine transformations, as well as the addition of outliers.

4.4. Benchmark Data

In addition to the simulations, RF, RerF, and RotRF were evaluated on 117 of the 121 datasets as described in (Fernandez-Delgado et al., 2014). The four remaining datasets were not used because their high dimensionality and large number of data points rendered the classifiers both time and space costly, particularly for RotRF. As in the previous section, transformations, with the exception of outliers, were applied to the datasets to observe their affects on performance of the three algorithms. Classifiers were trained on the entire training sets provided. For each data set, misclassification rates were again estimated by out of bag error. The number of trees used in each algorithm was 1000 for datasets having at most 1000 data points and 500 for datasets having greater than 1000 data points. All algorithms were trained using five different values for d , and error rates for each algorithm were selected as the minimum given by the five. For each dataset, relative performance ratios for each algorithm were computed by dividing the error rate of each algorithm by the minimum error rate of the three. The empirical cumulative distribution func-

tions of relative performance ratios for each algorithm were computed and plotted in Figure 4. Such plots are called performance profiles (Dolan & Moré, 2008). Performance profiles are useful in visualizing how frequently a particular algorithm wins, and when it loses, how frequently it loses by a certain amount. Areas under the curves (AUCs) were computed for each performance profile. This metric is valuable in that it doesn't necessarily indicate how frequently a particular algorithm is the best. Rather, it is an indicator of robustness. For instance, it is possible that an algorithm that wins on the majority of datasets can have a lower AUC than the other algorithms if, when it loses, it frequently loses by a large margin. On the original benchmarks (Figure 4 panel A), RerF achieves the largest AUC, followed by RF and lastly RotRF. The same trend is seen when scaling (panel C) and affine transformations (panel D) are applied. When the datasets are rotated, RotRF achieves the largest AUC, followed by RerF, and lastly RF.

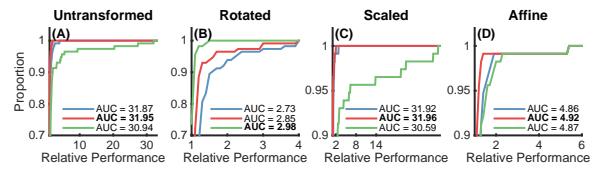


Figure 4. Performance profiles on benchmarks with various transformations applied (see section 4.4 for details). AUC denotes the area under the curve. RerF outperforms RF in all cases and RotRF in all cases except for rotation

4.5. Rank Transforming the Data and Fast Supervised Projections

When dimensions of X are linearly combined, we lose *scale and unit invariance*. We therefore note that random forests have a special property: they are invariant to monotonic transformations of the data applied to each coordinate in the ambient (observed) space. They are effectively operating on the order statistics, rather than actual magnitudes of coordinates. In other words, if we convert for each dimension the values of the samples to their corresponding ranks, random forests yield the exact same result. Therefore, for our second trick, we adopt the same policy, and “**pass to ranks**”, or rank transform, prior to doing anything else. We call RerFs that pass to ranks RerF(r).

Additionally, there is evidence that *using the data to construct A at each node* can significantly improve performance (Heath et al., 1993). However, previously proposed approaches for using the data require methods that are time- and space-intensive compared to standard random forests. We propose a simple strategy: “**compute the mean dif-**

ference vectors". In other words, given a two-class problem, let $\hat{\delta} = \hat{\mu}_0 - \hat{\mu}_1$, where $\hat{\mu}_c$ is the estimated class conditional mean. Under a spherically symmetric class conditional distribution assumption, $\hat{\delta}$ is the optimal projection vector. When there are $C > 2$ classes, the set of all pairwise distances between $\hat{\mu}_c$ and $\hat{\mu}_{c'}$ is of rank $C - 1$. Because RerFs are approximately rotationally invariant, we can simply compute all the class conditional means, and subtract each from one of them (we chose the $\hat{\mu}_c$ for which n_c , the number of samples in that class, is the largest). Thus, we construct \tilde{X} by concatenating A^\top with $\Delta = (\delta_{(2)} - \delta_{(1)}, \dots, \delta_{(C)} - \delta_{(1)})$, where $\delta_{(j)}$ is the δ vector for the class with the j^{th} largest number of samples, to obtain \tilde{A}^\top . Then we compute $\tilde{X} = \tilde{A}^\top \tilde{X}$. Computing this matrix is extremely fast, because it does not rely on costly singular value decompositions, matrix inversions, or convex problems. Thus, it nicely balances using the data to find good vectors, but not using much computational space or time. Furthermore, it also provides a set of dense projections to supplement the sparse projections of RerF. Denote RerFs that include this matrix RerF(d), and if they pass to ranks first, RerF(d+r).

Figure 5 shows classification error rates of RerF, RotRF, RerF(d), and RerF(d+r) against those of RF for each of the benchmark datasets in panels A-D, respectively. Panel A shows that RerF performs as well as RF for nearly all of the datasets, and does substantially better on a few. Panel B shows that RotRF does marginally worse in many cases, substantially worse in a few cases, and substantially better in a few cases. Panel C shows that RerF(d) performs marginally worse than RF on a moderate number of datasets, marginally worse on a moderate number, substantially worse on a couple of datasets, and substantially better on just one. Panel D shows that RerF(d+r) shows a similar trend to RerF(d+r). While using the mean difference vectors hurts the overall performance of RerF, it does not preclude the possibility of there being cases in which it is helpful. It is also worth experimenting with various ways of utilizing the mean difference vectors. For instance, we utilized them at every split node. However, only using them at random split nodes according to some scheme might have a different effect.

In addition to classification performance on the benchmark data, we were also interested in comparing training times of RF, RerF, and RotRF. In the same fashion as Figure 4, we plotted performance profiles for training times in Figure 6. Across all transformations, we see that RF and RerF perform comparably in terms of speed. Panels A and C show that RotRF performs the worst on the original benchmark data and the randomly scaled benchmark data. It performs well, however, on the rotated and affine transformed data. One possible explanation for this is that when the data is rotated, RF and RerF have a more difficult time finding good

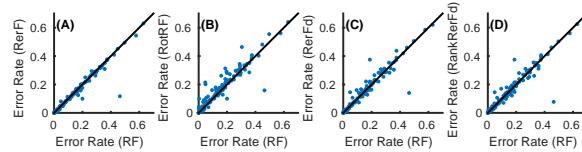


Figure 5. Scatter plots of the error rates for (A) RerF, (B) RotRF, (C) RerF(d), and (D) RerF(d+r) against that for RF on each benchmark dataset. The black line indicates the points of equal classification performance.

split directions, which could result in deeper trees (and subsequently more time spent in training).

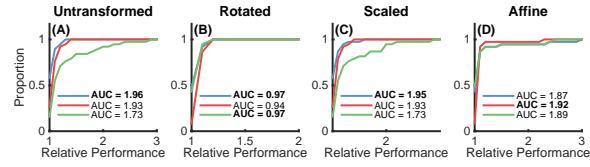


Figure 6. Performance profiles for classifier training times on benchmarks with various transformations applied. Performance of RF and RerF are comparable across all transformations. RotRF significantly underperforms RF and RerF on the untransformed and scaled data.

5. Conclusion and Future Work

We have proposed novel methods for constructing decision forests, which we call RerFs. We view these methods as special cases of a more general random projection forest, which include Breiman's original Forest-IC and Forest-RC, as well as previously proposed oblique decision forests. We have demonstrated in simulations that RerFs are especially well-suited for classification problems in which axis-parallel splits are suboptimal, and at the same time, have a large number of irrelevant features relative to relevant ones. This could explain RerF's excellent empirical performance on a suite of over 100 benchmark datasets, as real data often has the properties just described. Moreover, RerF, unlike other oblique decision forests, preserves the time and space complexity of Forest-IC, and is only marginally sensitive to scaling. Lastly, we propose two additional modifications. One involves supplementing sparse random projections with a particular form of fast supervised projections in the construction of candidate split directions. The other involves passing the data to ranks with the goal of making the classifier more robust to outliers.

Much work is still to be done with our proposed methods.

605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659

660
 661 The simplicity of RerFs suggest that they will be amenable
 662 to theoretical investigation, extending the work of Biau et
 663 al. (2008) to the RerF setting (?). Moreover, we hope
 664 that theoretical investigations will yield more insight into
 665 which distributions $f_A(\mathcal{D}_n)$ will be optimal under different
 666 distributional settings, both asymptotically and under finite
 667 sample assumptions. Even under the currently proposed
 668 $f_A(\mathcal{D}_n)$, our implementation has room for improvement.
 669 Although it has the same space and time complexity as RF,
 670 we will determine explicit constants, and improve our im-
 671 plementation accordingly. Indeed, our current implementa-
 672 tion is a proof-of-concept MATLAB implementation. We
 673 will utilize recent GPU and semi-external memory meth-
 674 ods to significantly speed up RerF (?). As with all decision
 675 forests, multiclass classification is but one exploitation task
 676 they can be used for; therefore, we will also extend this
 677 work to enable regression, density estimation, clustering,
 678 and hashing. We will provide open source code to enable
 679 more rigorous and extended analyses by the machine learn-
 680 ing community.

References

- 683 Blaser, R. and Fryzlewicz, P. Random rotation ensem-
 684 bles. 2015. URL [http://stats.lse.ac.uk/
 685 fryzlewicz/rre/rre.pdf](http://stats.lse.ac.uk/fryzlewicz/rre/rre.pdf).
- 686 Breiman, L. Random forests. *Machine Learning*, 4(1):5–
 687 32, October 2001.
- 688 Caruana, R., Karampatziakis, N., and Yessenalina, A. An
 689 empirical evaluation of supervised learning in high di-
 690 mensions. *Proceedings of the 25th International Con-
 691 ference on Machine Learning*, 2008.
- 692 Dolan, E. D. and Moré, J. J. Benchmarking optimization
 693 software with performance profiles. *Mathematical Pro-
 694 gramming*, 91:201–213, 2008.
- 695 Fernandez-Delgado, M., Cernadas, E., Barro, S., and
 696 Amorim, D. Do we need hundreds of classifiers to solve
 697 real world classification problems? *Journal of Machine
 698 Learning Research*, 15(1):3133–3181, October 2014.
- 699 Heath, D., Kasif, S., and Salzberg, S. Induction of oblique
 700 decision trees. *Journal of Artificial Intelligence Re-
 701 search*, 2(2):1–32, 1993.
- 702 Ho, T. K. The random subspace method for constructing
 703 decision forests. *Pattern Analysis and Machine Intelli-
 704 gence, IEEE Transactions on*, 20(8):832–844, Aug 1998.
 705 ISSN 0162-8828. doi: 10.1109/34.709601.
- 706 Li, P., Hastie, T. J., and Church, K. W. Very sparse random
 707 projections. In *Proceedings of the 12th ACM SIGKDD
 708 international conference on Knowledge discovery and
 709 data mining*, pp. 287–296. ACM, 2006.

- 710 Menze, B. H., Kelm, B.M., Splitthoff, D. N., Koethe,
 711 U., and Hamprecht, F. A. On oblique ran-
 712 dom forests. In Gunopulos, Dimitrios, Hofmann,
 713 Thomas, Malerba, Donato, and Vazirgiannis, Michalis
 714 (eds.), *Machine Learning and Knowledge Discovery
 715 in Databases*, volume 6912 of *Lecture Notes in Com-
 716 puter Science*, pp. 453–469. Springer Berlin Heidel-
 717 berg, 2011. ISBN 978-3-642-23782-9. doi: 10.1007/
 718 978-3-642-23783-6_29. URL http://dx.doi.org/10.1007/978-3-642-23783-6_29.
- 719 Rodriguez, J. J., Kuncheva, L. I., and Alonso, C. J. Rota-
 720 tion forest: A new classifier ensemble method. *Pattern
 721 Analysis and Machine Intelligence, IEEE Transac-
 722 tions on*, 28(10):1619–1630, 2006.
- 723 Tan, P. J. and Dowe, D. L. Mml inference of oblique de-
 724 cision trees. In *AI 2004: Advances in Artificial Intelli-
 725 gence*, pp. 1082–1088. Springer, 2005.
- 726 Trunk, G. V. A problem of dimensionality: A simple ex-
 727 ample. *Pattern Analysis and Machine Intelligence, IEEE
 728 Transactions on*, (3):306–307, 1979.