

Article

Enhancing Automatic Prediction of Spirality using SpArcFiRe's Spiral Arm Analysis and Random Forests

Pedro Silva ¹, Leon T. Cao ¹ and Wayne B. Hayes ^{1,*}

¹ Department of Computer Science, University of California, Irvine 92697-3435, USA; pedro.silva@uci.edu (P.S.); gcao@uci.edu (L.T.C.)

* Correspondence: whayes@uci.edu

Abstract: Automated machine classifications of galaxies are necessary because the size of upcoming surveys will overwhelm human volunteers. We improve upon existing machine classification methods by adding the output of SpArcFiRe to the inputs of a machine learning model. We use the human classifications from Galaxy Zoo 1 (GZ1) to train a random forest of decision trees to reproduce the human vote distributions of the Spiral class. We prefer the random forest model over other black box models like neural networks because it allows us to trace *post hoc* the precise reasoning behind the classification of each galaxy. We find that, across a sample of 470,000 Sloan galaxies that are large enough that details could be seen if they were there, the combination of SpArcFiRe outputs with existing SDSS features provides a better machine classification than either one alone on comparison to Galaxy Zoo 1. We suggest that adding SpArcFiRe outputs as features to any machine learning algorithm will likely improve its performance.

Keywords: galaxy morphology, machine learning; data analysis; object classification

1. Introduction

1.1. Motivation

The Hubble Ultra Deep Field (HUDF) represents about 1/13,000,000 of the celestial sphere and contains about 10,000 galaxies at least 100 of which have visible structure (by our own estimate), suggesting that the entire sky contains upwards of 10^9 galaxies with visible structure at the resolution and depth of the HUDF. To classify and quantitatively understand this number of galaxies will require automated methods.

SpArcFiRe¹ [1–3] is an algorithm designed to automatically extract structural information from the images of spiral galaxies. It was tested around a sample of 29,250 spiral galaxies from the Sloan Digital Sky Survey (SDSS), as selected by one of the PIs of the Galaxy Zoo project². The selection criteria were: $(GZ1_{P_S} + GZ1_{P_Z}) > 0.8$ OR $(GZ2_{FeaturesOrDisk} > 0.7$ AND $GZ2_{NotEdgeOn} > 0.7$ AND $GZ2_{spiral} > 0.8)$, where P_S is the fraction of human votes for S-wise (counterclockwise) spiral, P_Z is the fraction of human votes for Z-wise (clockwise) spiral, and *spiral* is the addition of P_S and P_Z , for each object, in either Galaxy Zoo 1 (GZ1) or Galaxy Zoo 2 (GZ2). This sample used the same magnitude limit as GZ1 (17.7 in the red band).

Even though some galaxy images (eg., elliptical galaxies or low-resolution spirals) do not have visible arms, we do not know in advance which images exhibit arms. For this reason, we run SpArcFiRe on *every* galaxy image, and our goal is to figure out when SpArcFiRe's output is meaningful, preferably using the output of SpArcFiRe itself. SpArcFiRe's job is to find spiral arms in spiral galaxies; often it also marks noise as spiral structure. Thus, we wish to recognize when a galaxy image has visible spiral

¹ SpArcFiRe stands for **S**Piral **A**RC **F**inder, and **R**Eporter and it is available at <http://sparcfire.ics.uci.edu>.

² Stephen Bamford, Personal Communication.

structure. Although ultimately we hope to develop an objective, quantitative, continuous measure of galaxy morphology, for now, we focus on the simple task of reproducing what we call the *spirality* of a galaxy image: from the GZ1 catalog [4,5], we define the *spirality* to be $P_{SP} = (GZ1_{P_S} + GZ1_{P_Z})$, representing the probability that there is any spiral structure visible for each object. We emphasize that spirality is a measure of the *image*, not the object. We are not trying to classify galaxies; we are trying to discern if a particular image exhibits spiral structure that is unlikely to be caused by noise. For example, although elliptical galaxies should be assigned a spirality of zero, an edge-on disk *also* should be assigned a spirality of zero, because spiral structure is not visible; thus, we wish to detect in both cases that SpArcFiRe’s output should not be interpreted as representing spiral structure.

Since humans introduce certain types of biases into the classification scheme (for example the chirality bias [5–7]), we also wish to “dilute” such biases even though we train our method on human classifications. We do this by carefully choosing which inputs we allow our code to use. For example, we allow SpArcFiRe’s measured pitch angle of spiral arms to be used as input to our machine learning classifier, but not the sign of the pitch angle [7], thus reducing chirality discrepancies to about 2 parts in 10,000 [3,7]. Our work follows up on existing work published in the astronomical literature [8–13].

1.2. Related Work

We compare against the most impactful and successful classifiers published in the astronomical literature: Banerji *et al.* [10] and Dieleman *et al.* [12]. The former was one of the first to apply Machine Learning to try and reproduce the human classifications of the GZ1 catalog [4] and the latter focuses specifically on reproducing the vote distribution of the GZ2 catalog [14], a regression problem, exactly like the approach we explore on this paper. The main difference here is that they are not concerned with the bias present in the dataset, so the smaller their Root Mean Squared Error (RMSE) is, “the better” their results are. In the recent Galaxy Zoo dataset releases, there has been an increased effort to eliminate human biases, but Hayes *et al.* [7] have proven that these datasets, in particular, GZ1, still contain biases so there is a trade-off between lowering the RMSE of a model and avoiding the introduction of such biases on the prediction.

Table 1. Non-SpArcFiRe input parameters we used, identical to those used in Banerji *et al.* [10], except for the absolute Magnitudes that also come from SDSS.

Name	Description
C&P set	colors and profile fitting
$dered_g - dered_r$	$(g - r)$ color, dereddened
$dered_r - dered_i$	$(r - i)$ color, dereddened
$deVAB_i$	de Vaucouleurs fit axial ratio
$expAB_i$	exponential fit axial ratio
$lnLexp_i$	exponential disk fit log likelihood
$lnLdeV_i$	de Vaucouleurs fit log likelihood
$lnLstar_i$	Star log likelihood
$absMag_X$	Absolute magnitudes in the 5 bands
AM set	adaptive moments
$petroR90_i / petroR50_i$	concentration
$mRrCc_i$	adaptive (+) shape measure
aE_i	adaptive ellipticity
$mCr4_i$	adaptive 4th moment
$texture_i$	texture parameter

Banerji *et al.* [10] present good results using neural networks. They classified Sloan galaxies in one of three categories: spiral, elliptical, and point sources/artifacts, using a neural network with inputs listed in Table 1. They found that on the entire sample of about 900,000 Sloan galaxies, they could reproduce the human GZ1 classifications in 92% of cases. Across a sample of brighter galaxies ($r < 17$), they correctly classify about 94% of galaxies. They do even better for a sample called the

“Gold sample”, in which galaxies are only included if the humans are themselves more than 80% confident in the classification. We do not believe the Gold sample comparison is meaningful, however, because it is crucial to know how good the machine learning classifier is when it *thinks* it is confident but is, in fact, mistaken, and the Gold sample completely disregards this aspect.³

Kaggle.com, a website devoted to machine learning competitions, offered £10,000 (GBP) to the algorithm which best minimized the RMSE between the automatic classification scheme and the human vote distribution for *Galaxy Zoo 2*. The winning entry was a Deep Learning algorithm using convolutional Neural Networks [12]. It had an RMSE of about 0.07 relative to the human GZ2 vote distribution. Although this result is closer to the human votes than our result presented below, we are concerned about the professional use of deep learning techniques for several reasons:

- We do not understand exactly what they are doing or how they are doing it, and research to better understand this aspect is still in its infancy [15,16]. Although we have some control over neural networks, we cannot learn from what they have learned, or learn from how they make their decisions, because a neural net is a near-complete “black box”.
- We would prefer an objective, quantitative system, with parameters that are understood and can be modified by professional astronomers, and decision trees seem better suited to this task.
- Decision trees are often used to measure the quality of features⁴ used to make a decision and thus are more suitable for our goals in this paper. This is not the case for Deep Neural Networks, which do not yet easily provide a similar measure for the features it used.

For these reasons, we prefer a method that can be understood, dissected, and whose individual decisions can also be understood and dissected, if necessary. Understanding these decisions can teach us about galaxy characteristics and morphology in ways that “black box” machine learning classifiers cannot. Figure 1 provides an example, giving a quantitative flavor of how color and magnitude provide information — both separately and together — about separating spirals and ellipticals, in ways a Neural Network cannot.

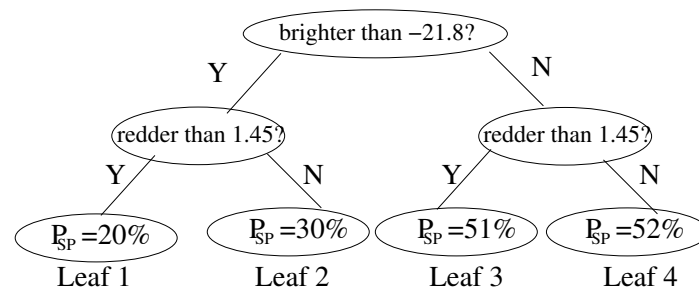


Figure 1. A simple 2-level, 2-parameter decision tree. The value of P_{SP} in each leaf node is the average probability of the training-set galaxies in that node being spiral, as measured by the GZ1 vote fraction. That value is then used as the predicted P_{SP} for galaxies in the non-training set that land in that node. With proper optimization (beyond the scope of this paper), larger trees with more features can produce more accurate predictions of P_{SP} , P_{EL} , etc.

Other interesting works published recently include Abd Elfattah *et al.* [18], which uses Neural Networks and Empirical Mode Decomposition to perform galaxy classification but uses a very small test set of 108 objects so it is hard to predict how their models would fare when trying to classify a

³ In essence, comparing against the Gold sample says “look how well we do when the humans pre-select the easy ones for us!” More formally, it disregards false positives — galaxies which the prediction is confident but is actually way off.

⁴ It is important to clarify that throughout this paper we will be using the term feature(s) to describe an individual measurable property or characteristic of a phenomenon being observed[17] as it is commonly done in the machine learning literature as opposed to features as seen in an image - like globular clusters.

much larger set of objects, like our test set. Kuminski *et al.* [11] makes a case for using “high-quality data” but we believe this will have the same issues as Banerji *et al.* [10]’s use of a “Gold Sample”. Applebaum and Zhang [19] uses an ensemble of Support Vector Machines to classify GZ2 galaxies achieving good results. Ferrari *et al.* [13] uses Linear Discriminant Analysis to classify galaxies from a couple different surveys.

All of the aforementioned work present good accuracies (\geq than 90%) but, except for Dieleman *et al.* [12], they are tackling this problem from a different point of view: they are all performing classification rather than regression. There have been tremendous advances in Machine Learning towards improving classifiers, and most of these papers make use of those techniques, but that is not the goal of our work. Whereas in classification one is concerned in finding a line that best separates two or more classes (in our case, spiral and non-spiral galaxies), in regression we seek to learn about the underlying distribution, in this case, how to put a probability on a galaxy being spiral⁵, and at present the GZ votes are the best way to do that. Usually more information is gleaned from a continuous distribution than a discrete classification — in particular a user of the output can choose a confidence threshold themselves for classification that is more suitable for a certain task rather than relying on the table creator’s subjective determination of where that threshold should lie. Peng *et al.* [20], for example, used regression for a task where they needed to analyze how spirality prediction degraded as a function of redshift, a task for which classification gives limited information.

For the sake of comparison we can turn our regressor into a classifier by choosing a boundary for the decision. If we choose that boundary to be 0.5, we will make our decision based on the majority vote, which mimics the choice of the Galaxy Zoo researchers in some releases [4]. That would give our regressor an accuracy of approximately 93% based on the test set presented in Table 5.⁶

2. Methods

We are mostly concerned with correctly predicting spirality (the probability of an image of a galaxy having visible spiral structure) for images of galaxies, in which spiral structure is visible, that have a reasonably high resolution. In particular, since SpArcFiRe is designed to discern spiral structure in disk galaxies, we are most interested in isolating disk galaxies in which spiral structure is visible. By a judicious eyeball study of images at the low end of resolution, we have subjectively determined that spiral structure is invisible in Sloan galaxies if the full major axis of the observable image is less than about 13 pixels, so we ignore any galaxy smaller than this. This is similar to the cutoff of 4.5 arcseconds petrosian radius used by the GZ1 team for galaxies with visible structure [4]. Also following GZ1, we cut off galaxies dimmer than magnitude 17.7 in the R band. This leaves about 470,000 Sloan galaxies.

We created models using Weka [21] which provides many machine learning algorithms, an easy-to-use interface, and the ability to create sophisticated standalone command-line classifiers once the model has been trained. Weka provides, among many algorithms, a Neural Network algorithm, and a Random Forest algorithm. Neural Networks have been used with success in similar tasks like the convolutional model used by Dieleman *et al.* [12]. These models excel in tasks where the input is spatially or temporal correlated like images or audio, so we briefly used the Neural Network algorithm to roughly reproduce the results of [10], having downloaded the same data they used from the Galaxy Zoo 1 survey [4,5], which was a treated sample of the Sloan Digital Sky Survey Data Release 6 (SDSS DR6) [22]. Since our machine learning algorithm uses the data only after SpArcFiRe has processed it, we found that Weka’s Random Forest model had a lower RMSE, and as described in the previous

⁵ High spirality is a strong indicator of a galaxy being spiral, but it’s not a *necessary* condition. Galaxies with low spirality may be edge-on spirals, ellipticals, low-resolution spirals, or even disk galaxies without spiral structure, such as the Sombrero Galaxy.

⁶ A higher accuracy can be achieved if we use a boundary below 0.5. Note that since there were 6 choices in GZ1, any vote receiving more than 1/6 of the votes can be a winning vote; for example, a vote of 40% could be considered a classification if all the other choices had less than 40% of votes. It is also possible to get better accuracy, using the same features, if we build a classifier rather than a regressor, but that’s outside the scope of this paper.

section, a Random Forest model (described below) makes decisions that are easier for us to dissect and learn from. For the most advanced tasks, we recreated the same random forest models using Julia [23]; the results using Weka and Julia are virtually identical since the underlying mechanisms are the same.

To provide context, we explain the general idea of random forests. The “forest” part refers to a set of decision trees. Each decision tree has a set of input parameters. At each level of the tree, one asks if a particular parameter is in a specific range. For example, one level of the decision tree may ask if the galaxy has an absolute magnitude brighter than 18; another level may ask if it has a color redder than 0. The tree can be very deep, and once we arrive at a leaf node, we have a set of galaxies that satisfy an exact set of characteristics across the parameters that lie along the decision path to that node. The process of optimizing the decision tree is beyond the scope of this paper, but the *goal* is to optimize the leaf nodes to precisely define whatever output characteristic we are trying to reproduce. In our case, we are trying to reproduce the GZ1 human vote distribution. For example, one leaf may represent all galaxies where the human votes for (elliptical, spiral, other) are close to (0.80, 0.19, 0.01). This helps us determine what characteristics lead a decision tree to classify a galaxy as spiral, elliptical, or other.

The “random” part of a random forest refers to the fact that each decision tree’s input parameters are chosen *randomly* from a larger set of input parameters provided by the user. The number of parameters to use for each tree is itself an integer parameter (fixed, in our case), as is the number of trees to use. Each tree effectively constitutes an “expert” in galaxy classification using its chosen set of parameters, and the forest is then a “mixture of experts”, in which a voting mechanism is used to come up with the final classification. A mixture of experts generally results in a much better classification than a single tree trained on all parameters, because the signal of each expert reinforces all the others, while the noise of the experts tends to cancel each other out. ([24] provides an excellent introduction to this idea.)

Figure 1 is a simple example of a two-parameter decision tree. In this example, we will apply it only to galaxies that are clearly either spiral or elliptical. However, rather than a discrete classification, our goal is to provide just one number for each galaxy: the probability that it is a spiral galaxy. We use two familiar parameters: color and absolute magnitude. It is well known that elliptical galaxies tend to be both brighter and redder than spirals. Given a training set of galaxies that are truly either spiral or elliptical and given the colors and magnitudes of each, we perform the following set of operations to generate a 2-parameter decision tree:

- Compute the mean magnitudes M_s, M_e for spirals and ellipticals, respectively.
- Compute the mean colors C_s, C_e for spirals and ellipticals, respectively.
- Compute a threshold color T_C intended to separate spirals from ellipticals; we will simply use the midpoint $T_C = (C_s + C_e)/2$.
- Similarly compute a threshold magnitude $T_M = (M_s + M_e)/2$.
- Now for each galaxy, first ask which side of the threshold its color is on, and then ask which side of the threshold its magnitude is on.
- This bins each galaxy into one of four leaf nodes, as in Figure 1.

As we can see, the results are correlated with the correct answers but not strongly so: dim, blue-ish galaxies only have a slightly greater than 50% chance of being spiral, although it is true that bright, reddish galaxies are correctly measured as unlikely to be spiral. Table 2 re-iterates this fact in more detail, and provides an example of another pair of features that provide a better classification scheme, although still only about 75% “correct” in total.

Table 2. Classification results for two-level, two-feature trees like that in Figure 1. Columns p_i represent the average fraction P_{SP} , across galaxies in leaf node i , of GZ1 humans who voted that object to be a spiral galaxy, across the training set. This value is then the assigned P_{SP} for any non-training-set galaxy placed in this leaf node. *correctAll*: assuming $P_{SP} > 0.5$ represent a positive spiral classification, the percentage across all galaxies of correct classifications; *SPcapture*: the fraction of true spirals that are captured by this classification scheme. *SPcontam*: the fraction of galaxies classified as spiral that are incorrectly classified. **Top row**: exactly the tree of Figure 1. **Second Row**: a pair that arguably performs better because it has a higher total correct classification, primarily because it has far less contamination of non-spirals, even though it has a smaller capture fraction. It demonstrates that we can have 75% correct classifications even with just a two-parameter, two-level tree. See Table 1 for the meaning of the input variables.

pair	p1	p2	p3	p4	correctAll	SPcapture	SPcontam
<i>rest_{ug}, MI</i>	0.20;	0.30;	0.52;	0.51	65.7%	42.0%	48.7%
<i>deVAB_i, MRrCc_i</i>	0.12;	0.65;	0.38;	0.85	74.3%	32.4%	14.8%

We begin to see significantly better results when we start to add features and levels in an individual tree. Table 3 lists the extra features, both from SpArcFiRe and elsewhere, that we use. Table 4 demonstrates how much better the classification gets as we increase the number of features and number of trees in the forest.

Table 3. Outputs from SpArcFiRe that are used as input features for our model, in addition to those from Table 1. See Davis and Hayes [2] for full descriptions of these parameters. Parameters labeled “DCO” are measured only across arcs of “dominant chirality only”—that is, arcs of the “wrong” chirality, which are likely to be noise, are not included. The parameter “arcLenAt50%” means: lay arcs end-to-end sorted longest to shortest, resulting in a line of total length L , and measure the length of the arc that lies at the point $L/2$ along the line. If the arms are short at $L/2$, then short arcs tend to suggest the galaxy is either flocculent or non-spiral, whereas a long arc at this point suggests a more grand-design spiral. The “rankAt50%” feature is similar, except this is the integer rank of the arc touching the $L/2$ point. If the ratio $((\text{diskAxisRatio}) / (\text{bulgeAxisRatio}))$ is close to 1, it is suggestive of an elliptical galaxy, whereas if this ratio is significantly less than one it suggests a spiral galaxy (since the bulge axis ratio tends to be 1 from any vantage point, but not so for the disk.)

Feature	Description
bar_scores	SpArcFiRe’s various bar detection scores
avg(abs(pa))-abs(avg(pa))	pitch angle-weighted chirality consistency across arms
numArcs > L	SpArcFiRe’s count of arms of various lengths
numDcoArcs > L	SpArcFiRe’s count of dominant-chirality-only arms of various lengths (see text)
totalNumArcs	total number of arcs found by SpArcFiRe
totalArcLen	total length of all arcs found by SpArcFiRe
avgArcLen	average arc length across arcs found by SpArcFiRe
arcLenAtNN%	length of arc at NN=25%, 50%, and 75% of total length of arcs (see text)
rankAtNN%	arc rank at NN=25%, 50%, and 75% of total length of arcs (see text)
bulgeAxisRatio	axis ratio of bulge, if present
diskAxisRatio	axis ratio of entire galaxy image; values $\lesssim 0.2$ suggest an edge-on spiral rather than elliptical
disk/bulgeRatio	disk to bulge ratio
diskBulgeAxisRatio	ratio of (diskAxisRatio) / (bulgeAxisRatio)
gaussLogLik	Gauss Log Likelihood of ellipse fit
likelihoodCtr	likelihood of the center of the ellipse fit
abs(pa_alen_avg)	average pitch angle of arms, length-weighted
abs(pa_alen_avg_DCO)	average pitch angle only of arms of dominant chirality
twoLongestAgree	chirality agreement of two longest arcs

We now explore in depth how many total trees should be in the forest, and how many randomly chosen features should be in each tree. Recall that the *total* number of features is fixed (and is 101 in

Table 4. Illustration of how the results of the classification improve as we allow more complex trees, and larger forests.

Total #features	features / tree	# of trees	actual feature	correct
1	1	1	Color only	65%
1	1	1	Magnitude only	65%
2	2	1	color + mag	75%
3	3	1	col,mag,arcs	85%
7	7	1	various	~90%
35	7	10	various	~95%
35	7	50	various	~97%
101	7	100	various	99.9%

our case), but that each decision tree chooses some random set of features. We will look at how both of these parameters change the results.

Presumably, the more features a particular tree uses, the better that tree will be, although more care needs to be put into training these models to avoid overfitting. Figure 2 plots the Pearson correlation between the GZ1 human vote proportion for P_{SP} , and our reproduction of that proportion, as a function of how many features are used by each tree. As can be seen, increasing the number of features used by each tree generally results in improvement. However, since each tree chooses a *random* subset of features, there is a bit of noise in the curve. It becomes less obvious that there is an improvement beyond about 35 features per tree, so we use 35 in our final results below. We also see that the entire curve moves up as the number of trees in the forest increases.

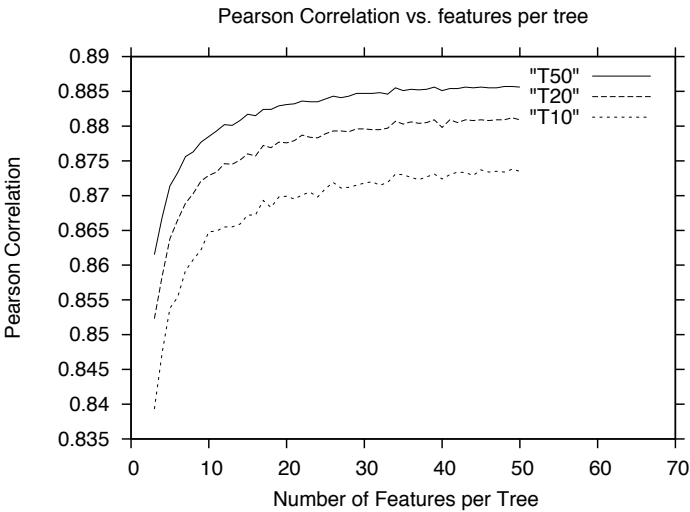


Figure 2. The Pearson correlation between the fraction of GZ1 humans voting for spiral, and our reproduction of that vote fraction, as a function of the number of features per tree that are chosen at random from the entire feature set. The three curves correspond to the cases where the total number of trees is 10, 20, or 50.

Similarly, we would expect that as the number of trees in the forest is increased, the result would get better. Essentially, as more “experts” weigh into the decision, the better the results should be. Figure 3 demonstrates that this is indeed the case. Furthermore, unlike the case of choosing features, the curve is pretty much monotonically increasing: it seems that more trees are always better [24]. In our results below, we use 150 total trees, each using 35 features out of our total set of 101 combined features from SpArcFiRe and SDSS.

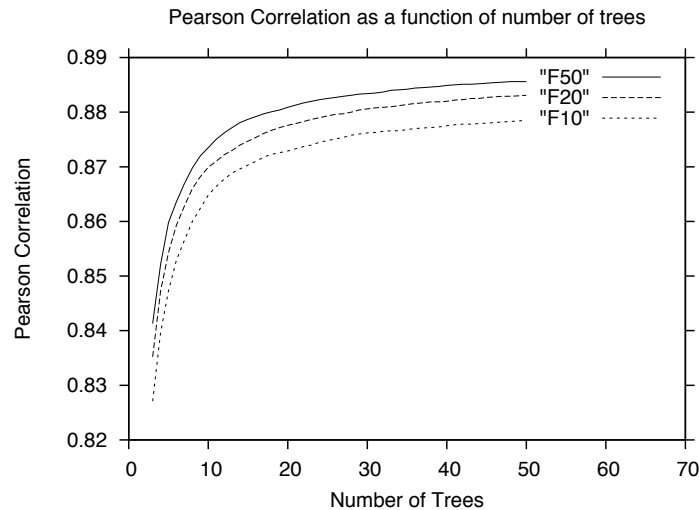


Figure 3. Similar to Figure 2, the Pearson correlation between the fraction of GZ1 humans voting for spiral, and our reproduction of that fraction, as a function of the number of trees in the forest. The three curves also show how the results change when the number of features per tree is 10, 20, or 50.

The advantage of this method over other more opaque methods such as Neural Networks, or SVM, is that once we get to a leaf node of the decision tree, we know *exactly* why each galaxy is in that node—we can follow the decisions down the tree and build a boolean expression that describes all the galaxies at that node. If we wish, we can then ask ourselves if the decision path makes sense; we can look at the galaxies at that node, and ask if they form an interesting set. This kind of detailed, explicit decision-making analysis is (currently) absent in other machine learning methods although very recent work has begun to study this question [15,16], and is what allows us to be more confident that biases are unlikely to creep into the classification scheme.

3. Results

As stated before, our goal is to test if adding SpArcFiRe's features to the set of input features will improve our ability to reproduce the vote distribution of GZ1 for spiral galaxies, so instead of classification, we are using regression to achieve our results. This means that rather than having a galaxy falling under a class (spiral, elliptical, and other) our output is the probability of an image of a galaxy having spiral structure. This value, between 0 and 1, is represented by the percentage of humans that agree that a certain galaxy has visible spiral structure. We represent this idea by making the sum of GZ1 values $P_S + P_Z$ as our *target variable*, and this is what we train our machine to reproduce — while simultaneously striving to eliminate the known P_S bias [6,7].

3.1. Measuring the quality of SpArcFiRe features

In the era of big data, machine learning scientists tend to agree that more is always better [24] but for some cases, this is not always true. Just adding features to a model does not guarantee that it will get better. Additional features might represent redundant information, which would not translate into more accurate classifiers for certain machine learning models, or worse, they would contribute to the curse of dimensionality [25]. In order to make sure we are adding meaningful information we further analyzed our features.

We built three different random forest models using the same hyperparameters (150 total trees, each using 35 features) but with different feature sets. Model 1 used only SDSS features, Model 2 used only SpArcFiRe features, and Model 3 used both sets of features (this is the model we discuss

Table 5. Predictions Confusion Matrix. The rows represent the number of objects that have a GZ1 spirality between a specific interval. The columns represent how many of those our Random Forest predicted in the same and different intervals. Notice that these numbers are only for the test set, thus a total of 45802 objects, which represent a more accurate measure of how our Random Forest would perform in real-world situations.

$P_{SP} \backslash RF_{SP}$	0.0-0.1	0.1-0.2	0.2-0.3	0.3-0.4	0.4-0.5	0.5-0.6	0.6-0.7	0.7-0.8	0.8-0.9	0.9-1.0	TOTAL
0.0-0.1	21238	5042	1512	522	169	45	17	1	1	1	28548
0.1-0.2	2471	1803	1145	594	254	111	28	5	0	0	6411
0.2-0.3	509	761	668	522	233	96	42	14	7	0	2852
0.3-0.4	184	345	424	348	209	120	58	23	5	2	1718
0.4-0.5	62	187	243	268	191	93	56	33	8	1	1142
0.5-0.6	47	128	215	200	188	145	76	38	11	3	1051
0.6-0.7	30	99	149	175	138	113	76	53	19	6	858
0.7-0.8	23	70	91	123	120	121	108	80	42	7	785
0.8-0.9	21	38	89	107	144	136	148	134	80	24	921
0.9-1.0	4	32	53	87	129	151	211	257	289	303	1516
TOTAL	24589	8505	4589	2946	1775	1131	820	638	462	347	45802

throughout the paper). We ran a 10-fold cross validation⁷ [26] in each one of those to get a more accurate measure of how those sets performed individually. Model 1 had a mean RMSE 0.1518, Model 2 had a mean RMSE of 0.1522, and Model 3 had a mean RMSE of 0.1404. For the tests and analysis made on this paper, we used the model with the lowest RMSE from the 10-fold cross validation used by model 3, which had an RMSE of 0.1374.

This demonstrates that SpArcFiRe features alone are just as good as SDSS features alone at predicting spirality. Furthermore, combining both sets has proven to increase the accuracy of our models. This is already an indication that there is valuable information in both feature sets.

Now let’s study our results in more detail. Table 5 shows our results, using both SDSS and SpArcFiRe’s features, for the test set in a 10x10 confusion matrix. Each row represents one of 10 bins holding galaxies in which a certain fraction of humans voted for that value of spirality; each column represents one of 10 identical bins containing the predicted spirality from our method. Thus, “correct” predictions (within 10% of the human vote) appear along the diagonal of the matrix. The first off-diagonal elements represent where our prediction was 10%-20% off, etc.; the far corners represent our worst predictions.

Notice that our model has high sensitivity and specificity rates, which means that when it predicts that an object is spiral or non-spiral with high confidence, the prediction is very likely correct. For example, let’s look at the case where our model predicts that an object is spiral with more than 90% of confidence, the penultimate column of the Table 5. If we consider a decision for spiral or non-spiral object being made above or below the 0.5 threshold, this gives us a sensitivity rate of more than 98%. The similar case happens for non-spiral predictions with more than 90% confidence (where $P_{SP} \leq 0.1$), the second column of the same table, in which, also considering a 0.5 threshold for a decision, our model gets more than 99% specificity rate.

In order to check which features seem to be the most important overall, we also created a feature ranking. As we have depicted in Figure 1, each node in a decision tree is a condition that splits the decision tree in two based upon a threshold in one variable. The measure used to make that decision is called impurity, and it is usually entropy for classification trees and variance for regression trees. It basically encodes how much information a particular feature, upon selection, adds to the decision

⁷ K-Fold cross validation is a method for measuring the quality of a learning algorithm by splitting the data into K buckets, training the algorithm in K-1 of these buckets and testing in the holdout bucket. We do this K times, each time holding out a different bucket and we report the average accuracy as the final accuracy of a model in that dataset.

Table 6. Top 10 best features for spirality prediction in decreasing order of importance. The standard deviation is measured across the 150 decision trees.

Feature	Score	Standard Deviation
Number of dominant-chirality-only arms equal or longer than 120	0.039	0.080
Absolute Magnitude in the Z band	0.031	0.020
De-reddened magnitude in the R band	0.029	0.019
De Vaucouleurs fit axial ratio i band	0.028	0.013
Number of dominant-chirality-only arms equal or longer than 85	0.022	0.061
Number of dominant-chirality-only arms equal or longer than 100	0.022	0.057
Number of arcs equal or longer than 120	0.022	0.057
Exponential fit axial ratio i band	0.021	0.009
De-reddened magnitude in the G band	0.021	0.015
Number of dominant-chirality-only arms equal or longer than 80	0.021	0.060

process. The more outputs a feature can separate, the higher its entropy is going to be, thus decreasing the impurity of the decision tree. So, we compute how much each feature decreases the weighted impurity of a tree. In our case, since we are using random forests, the impurity decrease from each feature can be averaged, and the features are ranked according to this measure [27].

Table 6 shows the top 10 features ranked by their importance along with the standard deviations of that score since this is an average over 150 decision trees. We can see that from the top 10 features 5 come from SDSS and 5 from SpArcFiRe, suggesting again that the two feature sets contribute roughly equally to the quality of the results. The 5 best SpArcFiRe features are all related to the number of arcs greater or equal to a certain amount of pixels, which is, not surprisingly, a strong indicative of the presence of spiral structure. Interestingly, in SpArcFiRe’s favor, the best feature overall is the number of dominant-chirality-only arms equal or longer than 120, which is 30% more relevant than the most relevant feature from SDSS.

Another way to visualize our results is to look at the Pearson correlation between our results and the GZ1 votes. Figure 4 shows this correlation represented in a graph where the x-axis represents the human votes and the y-axis our algorithm output, this time for all the 470,000 galaxies. Each red point represents one galaxy, and its (x,y) position represents our level of agreement. When x equals y, we are in complete agreement with the human votes. The clustering around the line $y = x$ suggests good agreement with GZ1. It is also notable that more than 98% of the galaxies have $|x - y| \leq 0.3$ and approximately 95% of the objects fall under $|x - y| \leq 0.2$.

In figure 5 we show some of our correctly classified objects. Those objects were cases where our model had a high agreement with the classifications provided by GZ1, and looking at the images we understand why. In 5a we display some of the spiral objects detected, while in 5b we show the non-spiral objects detected, which belong to the other classes of objects in GZ1: Elliptical, Merger, and Artefact, respectively.

It is important to understand what is going on in the 2% of objects that are outside of that scope. These objects are in the opposite corners of the off-diagonal in Table 5: 4 objects from the bottom left corner and 1 from the top right corner. These are that objects with a high disagreement: $|x - y| \geq 0.9$. From our total of 45802 galaxies in the test set, only 5 falls under this margin, and we show all of them in figure 6.

The top 4 rows depict the same problem: very faint arms that SpArcFiRe entirely failed to detect during the disk detection phase, so that it zoomed in past the arms, making it impossible for the arm detection code to find anything useful. This is a rare occurrence, and we are aware of this issue and are working on improving this specific step of the algorithm. The object on the bottom row is clearly a merger, and arm-like features are present, so our machine predicts a high spirality. One could argue that this is a *correct* prediction that the galaxy is not an elliptical galaxy, but the GZ1 humans correctly

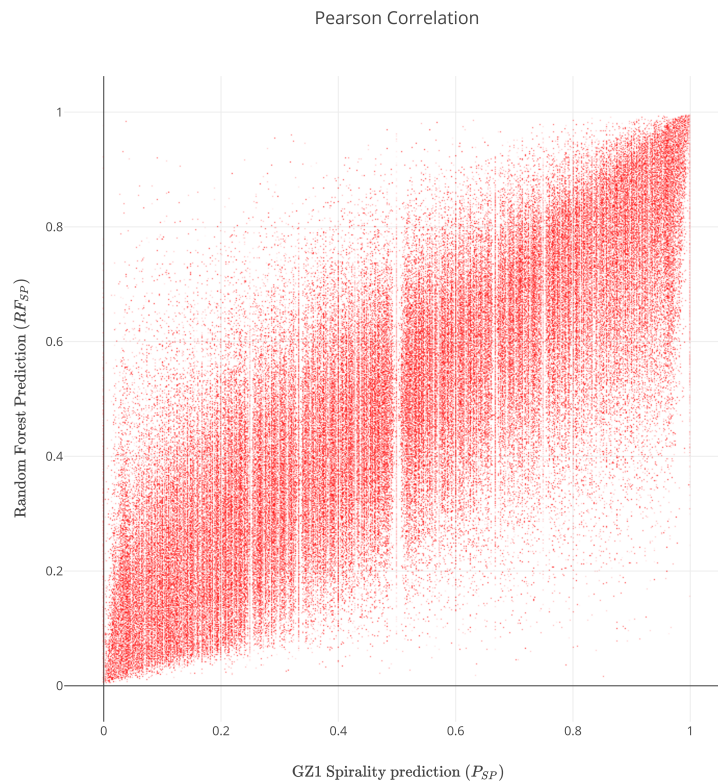


Figure 4. Scatter plot of our predicted spirality (vertical) vs. the fraction of GZ1 humans voting for spiral (horizontal). The Pearson correlation is 0.86, and the points cluster around the line $y = x$, depicting good agreement. Additionally, more than 98% of the galaxies have $|x - y| \leq 0.3$ and approximately 95% of the objects fall under $|x - y| \leq 0.2$. The vertical white lines appear because the fraction of human voters is a ratio of discrete integers.

marked it as a merger and thus not a spiral at all. Since our machine has not been trained to detect mergers, it is unclear whether this should count as a misclassification.⁸

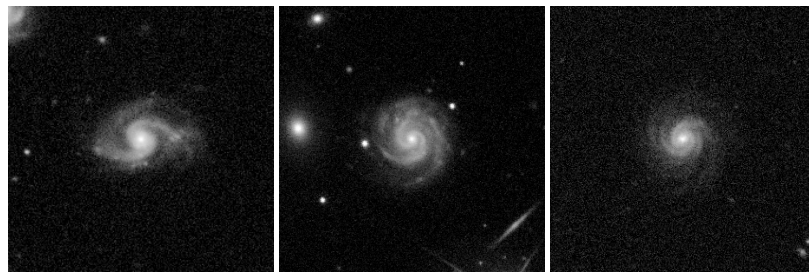
4. Conclusion

Our results show that it is possible to have a solid model that is in agreement with human classifiers above 90% of the time and also deal with the winding bias problem which was addressed in more detail in [7]. In this sense, we “filter” the errors made by humans while still retaining the useful knowledge provided by the Galaxy Zoo.

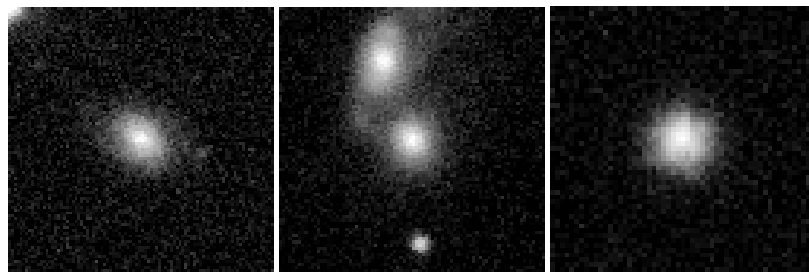
What differentiates this from previous work is the addition of SpArcFiRe’s output which adds more information to the objects we are discriminating and helps to decrease the amount of bias present in the classifications provided in GZ1. These results demonstrate that SpArcFiRe adds valuable information, rather than repeated, which can be used by automatic machine learning classifiers and regressors to achieve better results. We provided some insights on what these models find more descriptive for spiral galaxies demonstrating the most important parameters used by random forests in Table 6.

Further experimentation with SpArcFiRe information could contribute even more to automatic classification since this work focused on showing that its information could be useful when added to

⁸ One might argue that perhaps our “spirality” measure is more aptly called “non-ellipticity”.



(a) Spirals correctly detected by our model.



(b) Non-spirals correctly detected by our model.

Figure 5. Examples of images that had a high agreement of classification by both, our Random Forest Model and the GZ1 humans. (a) shows images of spirals where $P_{SP} \geq 0.90$ AND $|P_{SP} - F_{SP}| \leq 0.02$. Their SDSS IDs are, respectively, 1237654030325973054, 1237662306733916433, and 1237668298219847857. (b) shows images of non-spirals where $P_{SP} \leq 0.10$ AND $|P_{SP} - F_{SP}| \leq 0.02$. Their SDSS IDs are, respectively, 1237663529721397476, 1237661465447497940, and 1237661949201154382.

other data available from the surveys and SDSS. Although for the purposes of learning about what leads machine learning classifiers to differentiate galaxies as spiral or elliptical we focused on using random forests, further exploration with Neural Networks and Deep Learning have the potential to increase the final accuracy despite the fact that they would not add much in terms of knowledge for human classifiers. SpArcFiRe is also a project constantly evolving, with improvements in galaxy cropping and arc detection, which in turn reflects on how machine learning algorithms perform on this task.

Author Contributions: Conceptualization, P.S., L.T.C. and W.B.H.; Data curation, P.S. and L.T.C.; Formal analysis, P.S., L.T.C. and W.B.H.; Investigation, P.S. and L.T.C.; Methodology, P.S. and L.T.C.; Project administration, W.B.H.; Resources, W.B.H.; Software, P.S. and L.T.C.; Supervision, W.B.H.; Validation, P.S.; Visualization, P.S., L.T.C. and W.B.H.; Writing – original draft, P.S., L.T.C. and W.B.H.; Writing – review editing, P.S. and W.B.H..

Funding: This research received no external funding.

Acknowledgments: This work was supported by CAPES (Coordination for the Improvement of Higher Education Personnel - Brazil) through the Science Without Borders fellowship for Ph.D. Studies awarded to Pedro Silva.

Funding for the SDSS and SDSS-II has been provided by the Alfred P. Sloan Foundation, the Participating Institutions, the National Science Foundation, the U.S. Department of Energy, the National Aeronautics and Space Administration, the Japanese Monbukagakusho, the Max Planck Society, and the Higher Education Funding Council for England. The SDSS Web Site is <http://www.sdss.org/>.

The SDSS is managed by the Astrophysical Research Consortium for the Participating Institutions. The Participating Institutions are the American Museum of Natural History, Astrophysical Institute Potsdam, University of Basel, University of Cambridge, Case Western Reserve University, University of Chicago, Drexel University, Fermilab, the Institute for Advanced Study, the Japan Participation Group, Johns Hopkins University, the Joint Institute for Nuclear Astrophysics, the Kavli Institute for Particle Astrophysics and Cosmology, the Korean Scientist Group, the Chinese Academy of Sciences (LAMOST), Los Alamos National Laboratory, the Max-Planck-Institute for Astronomy (MPIA), the Max-Planck-Institute for Astrophysics (MPA), New Mexico State University, Ohio State University, University of Pittsburgh, University of Portsmouth, Princeton University, the United States Naval Observatory, and the University of Washington.

Conflicts of Interest: The authors declare no conflict of interest.

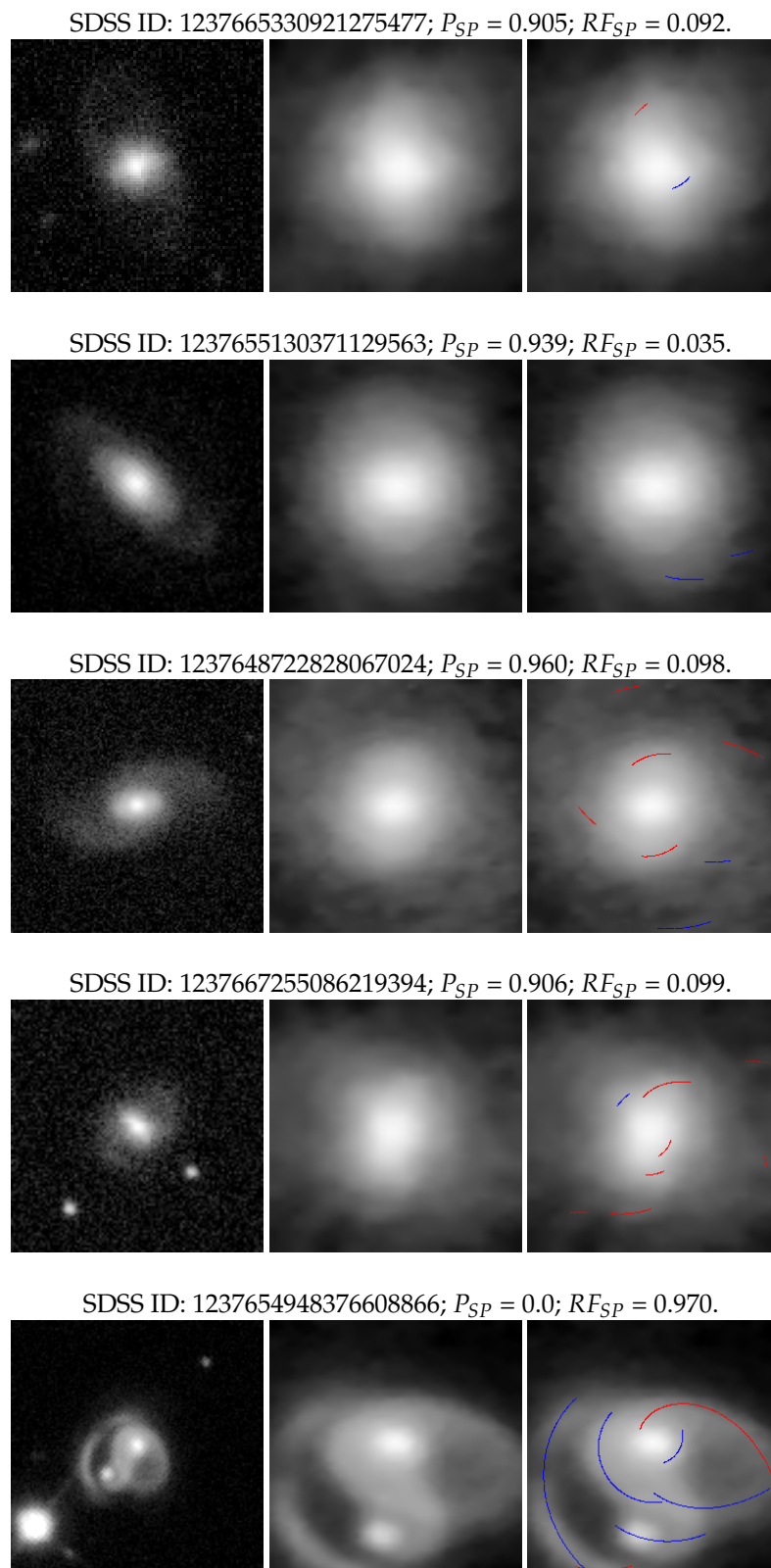


Figure 6. Grossly Misclassified Objects. In sets of 3, from left to right column, the images show the Original Input Image, the same image automatically cropped by SpArcFiRe, and the spiral Arcs detected on the image (if any). The SDSS Object IDs, the GZ1 Spirality prediction (P_{SP}), and our Random Forest Prediction (RF_{SP}) are shown above each trio of images. In all but the last, the problem is low-surface-brightness arms, which we know about and are working on this issue. Despite the disagreement in the 5th object, a merger, spiral structure is indeed present.

References

1. Davis, D.; Hayes, W. Automated quantitative description of spiral galaxy arm-segment structure. *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, 2012, pp. 1138–1145. doi:10.1109/CVPR.2012.6247794.
2. Davis, D.R.; Hayes, W.B. SpArcFiRe: Scalable Automated Detection of Spiral Galaxy Arm Segments. *The Astrophysical Journal* **2014**, *790*, 87.
3. Davis, D.R. Fast Approximate Quantification of Arbitrary Arm-Segment Structure in Spiral Galaxies. Phd thesis, University of California, Irvine, 2014.
4. Lintott, C.J.; Schawinski, K.; Slosar, A.; Land, K.; Bamford, S.P.; Thomas, D.; Raddick, M.J.; Nichol, R.C.; Szalay, A.; Andreescu, D.; Murray, P.; Vandenberg, J. Galaxy Zoo: morphologies derived from visual inspection of galaxies from the Sloan Digital Sky Survey. *Monthly Notices of the Royal Astronomical Society* **2008**, *389*, 1179–1189.
5. Lintott, C.J.; Schawinski, K.; Bamford, S.P.; Slosar, A.; Land, K.; Thomas, D.; Edmondson, E.; Masters, K.; Nichol, R.C.; Raddick, M.J.; Szalay, A.; Andreescu, D.; Murray, P.; Vandenberg, J. Galaxy Zoo 1: data release of morphological classifications for nearly 900 000 galaxies. *Neural Networks* **2011**, *410*, 166–178.
6. Land, K.; Slosar, A.; Lintott, C.J.; Andreescu, D.; Bamford, S.P.; Murray, P.; Nichol, R.; Raddick, M.J.; Schawinski, K.; Szalay, A.; Thomas, D.; Vandenberg, J. Galaxy Zoo: the large-scale spin statistics of spiral galaxies in the Sloan Digital Sky Survey. *Monthly Notices of the Royal Astronomical Society* **2008**, *388*, 1686–1692.
7. Hayes, W.B.; Davis, D.R.; Silva, P. On the nature and correction of the spurious S-wise spiral galaxy winding bias in Galaxy Zoo 1. *Monthly Notices of the Royal Astronomical Society* **2017**, *466*, 3928–3936.
8. Shamir, L. Automatic morphological classification of galaxy images. *Monthly Notices of the Royal Astronomical Society* **2009**, *399*, 1367–1372, [http://oup/backfile/content_public/journal/mnras/399/3/10.1111/j.1365-2966.2009.15366.x/1/mnras0399-1367.pdf]. doi:10.1111/j.1365-2966.2009.15366.x.
9. Huertas-Company, M.; Aguerri, J.A.L.; Bernardi, M.; Mei, S.; Sánchez Almeida, J. Revisiting the Hubble sequence in the SDSS DR7 spectroscopic sample: a publicly available Bayesian automated classification. *Astronomy & Astrophysics* **2010**, *525*, A157.
10. Banerji, M.; Lahav, O.; Lintott, C.J.; Abdalla, F.B.; Schawinski, K.; Bamford, S.P.; Andreescu, D.; Murray, P.; Raddick, M.J.; Slosar, A.; Szalay, A.; Thomas, D.; Vandenberg, J. Galaxy Zoo: Reproducing galaxy morphologies via machine learning. *Monthly Notices of the Royal Astronomical Society* **2010**, *406*, 342–353.
11. Kuminski, E.; George, J.; Wallin, J.; Shamir, L. Combining Human and Machine Learning for Morphological Analysis of Galaxy Images. *Publications of the Astronomical Society of the Pacific* **2014**, *126*, 959–967.
12. Dieleman, S.; Willett, K.W.; Dambre, J. Rotation-invariant convolutional neural networks for galaxy morphology prediction. *Monthly Notices of the Royal Astronomical Society* **2015**, *450*, 1441–1459.
13. Ferrari, F.; de Carvalho, R.R.; Trevisan, M. Morfometryka - A New Way of Establishing Morphological Classification of Galaxies. *Astrophysical Journal* **2015**, *814*, 55.
14. Willett, K.W.; Lintott, C.J.; Bamford, S.P.; Masters, K.L.; Simmons, B.D.; Casteels, K.R.V.; Edmondson, E.M.; Fortson, L.F.; Kaviraj, S.; Keel, W.C.; Melvin, T.; Nichol, R.C.; Raddick, M.J.; Schawinski, K.; Simpson, R.J.; Skibba, R.A.; Smith, A.M.; Thomas, D. Galaxy Zoo 2: detailed morphological classifications for 304 122 galaxies from the Sloan Digital Sky Survey. *mnras* **2013**, *435*, 2835–2860, [[arXiv:astro-ph/1308.3496](https://arxiv.org/abs/astro-ph/1308.3496)]. doi:10.1093/mnras/stt1458.
15. Ribeiro, M.T.; Singh, S.; Guestrin, C. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*; ACM: New York, NY, USA, 2016; KDD '16, pp. 1135–1144. doi:10.1145/2939672.2939778.
16. Nguyen, A.; Yosinski, J.; Clune, J. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. *STOC '97 Proceedings of the twenty-ninth annual ACM symposium on Theory of computing*. University of Wyoming, Laramie, United States, IEEE, 2015, pp. 427–436.
17. Bishop, C.M. *Pattern Recognition and Machine Learning (Information Science and Statistics)*, 1st ed. 20 ed.; Springer, 2007.

18. Abd Elfattah, M.; Elbendary, N.; Elminir, H.K.; Abu El-Soud, M.A.; Hassanien, A.E. Galaxies image classification using empirical mode decomposition and machine learning techniques. 2014 International Conference on Engineering and Technology (ICET. IEEE, 2014, pp. 1–5.
19. Applebaum, K.; Zhang, D. Classifying Galaxy Images through Support Vector Machines. 2015 IEEE International Conference on Information Reuse and Integration (IRI). IEEE, 2015, pp. 357–363.
20. Peng, T.; English, J.E.; Silva, P.; Davis, D.R.; Hayes, W.B. SpArcFiRe: morphological selection effects due to reduced visibility of tightly winding arms in distant spiral galaxies. *arXiv.org* **2017**, [1707.02021].
21. Hall, M.; Frank, E.; Holmes, G.; Pfahringer, B.; Reutemann, P.; Witten, I.H. The WEKA data mining software: an update. *ACM SIGKDD Explorations Newsletter* **2009**, *11*, 10–18.
22. Adelman-McCarthy, J.K.; Agüeros, M.A.; Allam, S.S.; Prieto, C.A.; Anderson, K.S.J.; Anderson, S.F.; Annis, J.; Bahcall, N.A.; Bailer-Jones, C.A.L.; Baldry, I.K.; Barentine, J.C.; Bassett, B.A.; Becker, A.C.; Beers, T.C.; Bell, E.F.; Berlind, A.A.; Bernardi, M.; Blanton, M.R.; Bochanski, J.J.; Boroski, W.N.; Brinchmann, J.; Brinkmann, J.; Brunner, R.J.; Budavári, T.; Carliles, S.; Carr, M.A.; Castander, F.J.; Cinabro, D.; Cool, R.J.; Covey, K.R.; Csabai, I.; Cunha, C.E.; Davenport, J.R.A.; Dilday, B.; Doi, M.; Eisenstein, D.J.; Evans, M.L.; Fan, X.; Finkbeiner, D.P.; Friedman, S.D.; Frieman, J.A.; Fukugita, M.; Gänsicke, B.T.; Gates, E.; Gillespie, B.; Glazebrook, K.; Gray, J.; Grebel, E.K.; Gunn, J.E.; Gurbani, V.K.; Hall, P.B.; Harding, P.; Harvanek, M.; Hawley, S.L.; Hayes, J.; Heckman, T.M.; Hendry, J.S.; Hindsley, R.B.; Hirata, C.M.; Hogan, C.J.; Hogg, D.W.; Hyde, J.B.; Ichikawa, S.i.; Ivezić, Ž.; Jester, S.; Johnson, J.A.; Jorgensen, A.M.; Jurić, M.; Kent, S.M.; Kessler, R.; Kleinman, S.J.; Knapp, G.R.; Kron, R.G.; Krzesinski, J.; Kuropatkin, N.; Lamb, D.Q.; Lampeitl, H.; Lebedeva, S.; Lee, Y.S.; Leger, R.F.; Lépine, S.; Lima, M.; Lin, H.; Long, D.C.; Loomis, C.P.; Loveday, J.; Lupton, R.H.; Malanushenko, O.; Malanushenko, V.; Mandelbaum, R.; Margon, B.; Marriner, J.P.; Martínez-Delgado, D.; Matsubara, T.; McGehee, P.M.; McKay, T.A.; Meiksin, A.; Morrison, H.L.; Munn, J.A.; Nakajima, R.; Eric H Neilsen, J.; Newberg, H.J.; Nichol, R.C.; Nicinski, T.; Nieto-Santisteban, M.; Nitta, A.; Okamura, S.; Owen, R.; Oyaizu, H.; Padmanabhan, N.; Pan, K.; Park, C.; John Peoples, J.; Pier, J.R.; Pope, A.C.; Purger, N.; Raddick, M.J.; Fiorentin, P.R.; Richards, G.T.; Richmond, M.W.; Riess, A.G.; Rix, H.W.; Rockosi, C.M.; Sako, M.; Schlegel, D.J.; Schneider, D.P.; Schreiber, M.R.; Schwobe, A.D.; Seljak, U.; Sesar, B.; Sheldon, E.; Shimasaku, K.; Sivarani, T.; Smith, J.A.; Snedden, S.A.; Steinmetz, M.; Strauss, M.A.; SubbaRao, M.; Suto, Y.; Szalay, A.S.; Szapudi, I.; Szkody, P.; Tegmark, M.; Thakar, A.R.; Tremonti, C.A.; Tucker, D.L.; Uomoto, A.; Berk, D.E.V.; Vandenberg, J.; Vidrih, S.; Vogeley, M.S.; Voges, W.; Vogt, N.P.; Wadadekar, Y.; Weinberg, D.H.; West, A.A.; White, S.D.M.; Wilhite, B.C.; Yanny, B.; Yocum, D.R.; York, D.G.; Zehavi, I.; Zucker, D.B. The Sixth Data Release of the Sloan Digital Sky Survey. *The Astrophysical Journal Supplement Series* **2008**, *175*, 297–313.
23. Bezanson, J.; Edelman, A.; Karpinski, S.; Shah, V.B. Julia: A Fresh Approach to Numerical Computing. *SIAM Review* **2017**, *59*, 65–98, [http://dx.doi.org/10.1137/141000671]. doi:10.1137/141000671.
24. Sibley, C. More Is Always Better: The Power Of Simple Ensembles, 2012. <http://www.overkillanalytics.net/more-is-always-better-the-power-of-simple-ensembles/>.
25. Jensen, R.; Shen, Q. Are More Features Better? A Response to *Attributes Reduction Using Fuzzy Rough Sets*. *IEEE Transactions on Fuzzy Systems* **2009**, *17*, 1456–1458.
26. Refaeilzadeh, P.; Tang, L.; Liu, H. Cross-Validation. In *Encyclopedia of Database Systems*; Springer New York: New York, NY, 2016; pp. 1–7.
27. Saabas, A. Selecting good features Part III: random forests, 2014.