

Tree Based Ensembles for Predicting Survival from Thoracic Surgery

Samuel Jackson, Aberystwyth University

I. INTRODUCTION

Thoracic surgery is a major invasive surgery involving operating on the lungs of a patient. The authors of ref. [1] collected several pieces of possibly relevant data on a number of patients who went on to have thoracic surgery. The data also includes a record of whether a given patient survived for longer than one year after the surgery. This paper looks at using a reduced subset of the features and patients from the dataset in [1] to classify patients based on whether or not they will survive for one year after the surgery. This paper compares four different classifiers: Random Forests [2], Extremely Randomised Trees [3], AdaBoost [4], and Gradient Boosting [5].

The format of the rest of this paper is structured as follows: section II outlines the preprocessing steps performed on the dataset and describes the classifiers used. Section III presents the performance of the classifiers on the dataset and explores different feature sets. The effects of resampling the dataset are also shown. Section IV discusses the results and presents possible justification for the performance based on the properties of the classifier and dataset. Finally, a summary and of possible future directions are discussed in section V.

II. METHODS

A. Dataset and Preprocessing

The thoracic surgery dataset used consists of 16 predictors and 300 instances. Table I gives a description of each predictor derived from the original UCI dataset repository [6]. The dataset includes a mixture of both categorical (nominal and ordinal) and continuous data. The final (17^{th}) column of the dataset is the binary class label with value 0 if the patient survived and 1 if they died within one year of surgery.

Several initial observations can be made about dataset prior to any preprocessing steps. One key thing to note about the dataset as a whole is that there is a slight imbalance between the two classes. Only 28% of the dataset is of the positive class (28% of patients died). While this imbalance is not extreme, it can have repercussions for the performance of the classifiers. The accuracy paradox [7] states that a classifier with high accuracy can be built from highly imbalanced data by always predicting the negative class.

The predictor PRE32 is zero for all of the patients in the training dataset. This predictor therefore has zero variation and will not help to discriminate between instance. PRE32 is therefore discarded during preprocessing.

PRE5 appears to have some extreme values. PRE5 corresponds to the FEV1 measure. This would suggest that some

TABLE I
DESCRIPTION OF COLUMNS IN THE THORACIC SURGERY DATASET

| Column | Type | Description |
|--------|---------|--|
| DGN | Nominal | Diagnosis: Specific combination of ICD-10 codes for primary and secondary as well multiple tumours if any (DGN3, DGN2, DGN4, DGN6, DGN5, DGN8, DGN1) |
| PRE4 | Numeric | Forced vital capacity (FVC) |
| PRE5 | Numeric | Volume that has been exhaled at the end of the first second of forced expiration (FEV1) |
| PRE6 | Ordinal | Performance status - Zubrod scale (PRZ2, PRZ1, PRZ0) |
| PRE7 | Nominal | Pain before surgery (T,F) |
| PRE8 | Nominal | Haemoptysis before surgery (T,F) |
| PRE9 | Nominal | Dyspnoea before surgery (T,F) |
| PRE10 | Nominal | Cough before surgery (T,F) |
| PRE11 | Nominal | Weakness before surgery (T,F) |
| PRE14 | Ordinal | T in clinical TNM - size of the original tumour, from OC11 (smallest) to OC14 (largest) (OC11, OC14, OC12, OC13) |
| PRE17 | Nominal | Type 2 DM - diabetes mellitus (T,F) |
| PRE25 | Nominal | Peripheral arterial diseases (PAD) (T,F) |
| PRE30 | Nominal | Smoking (T,F) |
| PRE32 | Nominal | Asthma (T,F) |
| AGE | Numeric | Age at surgery |
| Risk1Y | Nominal | 1 year survival period - (T)true value if died (T,F) (Class Label) |

patients have an unusually high forced expiration volume. Also, all of the outliers are of the same class. This could cause the classifiers to fit to noise rather than to properly generalise. These instances were therefore removed from the dataset. No reduction in performance was witnessed during cross validation for all classifiers after their removal.

The feature DGN is a nominal categorical predictor. This feature was transformed into series of new features via one hot encoding. Each new predictor is a binary feature which is one if the patient falls into the category and zero otherwise. The original DGN feature is drop after the 7 new binary features are created.

Finally, after all preprocessing is complete, a Random Forest is trained on the dataset (with default parameters) and the resulting variable importance measure is computed. Any features with a variable importance of zero are dropped. The variable importance of the preprocessed features (before any are dropped) is shown in figure 1.

B. Classifiers

Four classifiers were chosen for use on the dataset. The four classifiers used are Random Forests, Gradient Boosting, AdaBoost, and Extra Trees. The implementations of all four classifiers are taken directly from the scikit-learn library

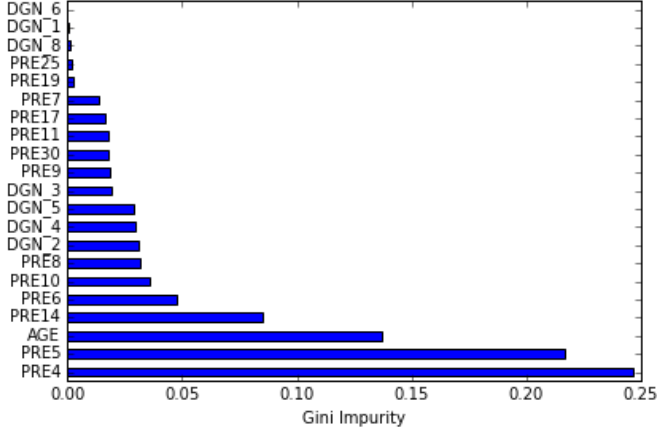


Fig. 1. Feature importance for all of the features after preprocessing. Variables with the prefix DGN_ are the one-hot encoded versions of the original DGN variable. The numerical predictors PRE4, PRE5, and AGE are shown to be the most important to survival prediction.

[8]. All of these methods are ensemble methods. Ensemble methods compose together multiple weak learners to produce a single strong learner. In this paper, all of the algorithms used are also tree based (although this is not necessarily the case for Gradient Boosting and AdaBoost). This means that for each classifier the base learner is a decision tree.

Random forests [2] are perhaps the simplest method of the four. Random forests are simply a collection of n decision trees which are individually trained on the data. The “random” in the name comes from the fact that each tree is trained on both a random sample of the dataset (tree bagging) and also on a random subset of the features (feature bagging). The decision tree weak learners typically overfit to the data. However because the results of all trees are averaged over the variance in the final model’s prediction is significantly reduced. The random element for both training instances and features is used to prevent many highly correlated trees from occurring and therefore reduce overfitting.

Extremely Randomised trees [3] (also called Extra-Trees) takes the randomisation aspect of random forests one step further. Random features are still used but training data is not bootstrapped. This aims to reduce model bias. Additionally a random split for each feature in the subset of features is chosen instead of the just computing the optimal feature and split combination. This has the benefit of potentially lowering the variance in the model.

The AdaBoost algorithm [4] is a generalised method for combining the performance of many weak learners. AdaBoost stands for “adaptive boosting” and boosting is a core component in the training stage. AdaBoost works by first fitting a weak learner to the training dataset and classifying each instance. The error in the classification can be used to re-weight each example in the dataset. Misclassified samples are therefore more likely to be classified correctly in future iterations. Likewise, instances that are correctly classified can be weighted much lower, as they are easier to correctly identify. AdaBoost can be seen as an additive method in that each tree is built on the error of the previous one. Adaboost

does not necessarily have to be used with decision trees, but in this paper only decision tree based AdaBoost is considered.

Gradient Boosting Machines [5] also use a boosting based approach to learning. Like AdaBoost they are an additive method that iteratively fits a collection of weak learners. Where the two differ is in the way that instances are “weighted”. The weighting function in AdaBoost can be seen as a special type of loss function. The gradient of a differentiable loss function can be used to steer the search towards the optimum decision function. In gradient boosting machines the new function added to the mix is the one which is the closest to parallel with the negative gradient of the observed data.

These algorithms were chosen to showcase a broad range of different ensemble algorithms. Ensemble methods often outperform a single strong learner due to the diversity present in the model. All of the models are also decision tree based which work well with a mixture of continuous and discrete values like those present in the dataset.

Two common techniques used in training ensemble algorithms are bagging and boosting. This paper compares two algorithms based on boosting (AdaBoost and Gradient Boosting) and one based on bagging (Random Forests). Extra-Trees also utilise bagging for feature subsets. These seem like good candidates based on the dataset for two reasons: 1) bagging can be used address the imbalance in the dataset by equally resampling each of the datasets and 2) there doesn’t seem to be a clear separation between classes in the dataset which potentially makes boosting a good technique as it should help to push the algorithm towards classifying the difficult missed examples.

C. Hyperparameters & Tuning

Before all experiments were carried out on the dataset, the hyper-parameters of each classifier were tuned to hopefully achieve optimum performance. The values and number hyper-parameters are dependant both on the implementation of the classifier and the dataset itself. If there is more than one hyper-parameter for a classifier (as is the case with all classifiers used here) then ideally combinations of all hyper-parameters should be explored. Sadly, this means that the space of potential hyper-parameters choices explodes as the number of hyper-parameters increases.

Due to the relatively small size of the dataset, the space of potential parameters for each classifier is explored using a grid search. In a grid search, a selection of hyper-parameter values are explicitly enumerated. Each potential value for a hyper-parameter is tested in combination with every other hyper-parameter value. The speed of the grid search is bearable due to the classifier being relatively quick to train on this small dataset. The performance of a set of parameters was evaluated using stratified k -fold cross validation with ROC AUC as the scoring metric.

For Random Forests a single grid search was performed over the tree parameters *max_depth*, *max_features*, *min_samples_split* and *min_samples_leaf*. *max_depth* and *max_features* were trained over the range 2 - 20 in steps

of 3. *min_samples_split* and *min_samples_leaf* were trained over all values in the range 1 - 5. The number of trees used was fixed to 50 during this search. This is because a small number of trees will be quick to train (and hence the search will complete faster). Generally speaking the performance of the forest should improve as the number of trees increases, so this can be trained afterwards.

The results of tuning the tree parameters showed that *min_samples_split* and *min_samples_leaf* should be set to 1. This seems logical due to the small number of positive samples. *max_depth*, *max_features* were optimised as 16 and 5 respectively. These seem reasonable given the low number of predictive features and the fact that trees in random forests should typically overfit (hence the large maximum depth). After this trial another grid search was performed to find the optimum number of trees over the range 50 - 500 in steps of 50. This suggested that 100 trees should be used.

The training procedure for Extremely Random Trees was identical to Random Forests with the results being similar. After tuning 200 trees were used and *max_features* was set to 16 and *max_depth* set to 19.

Adaboost was tuned by fixing the maximum depth of the decision tree and performing a grid search over the number of trees (50 - 1000 in steps of 50) and the learning rate (with values 0.1, 0.5, 0.01, and 0.005) together. This was repeated for multiple values of the maximum depth (tested with values 2, 4, and 6). After tuning 400 trees were used with a *learning_rate* of 0.5 and *max_depth* of 4.

Gradient Boosting has many parameters that need to be explored, many of which can interact with one another and tuning in the wrong order can lead to poor results. The number of parameters can also be awkward to train due to the speed of training. The tuning procedure for gradient boosting was therefore as follows:

- Fix all of the parameters to be reasonable initial guesses and fix the learning rate to be quite high (0.1).
- Find the optimum number of estimators for the given learning rate (searched over the range 20 - 150 in steps of 10).
- Tune the tree based parameters *max_depth* and *min_samples_split* (searched over the ranges 5-16 in steps of 2 and 1-20 in steps of 3 respectively).
- Tune *max_features* (5-20 in steps of 2)
- Tune the subsample ratio (values: 0.6, 0.7, 0.75, 0.8, 0.85, and 0.9).
- Finally using all previously tuned parameters increase the number of estimators while simultaneously decreasing the learning rate.

The final values for the tuned parameters used are shown in table II.

III. RESULTS

A. Performance Evaluation

Each classifier in section II-B was trained using stratified 5-fold cross validation. Stratification was performed to ensure that there was a representative sample of positive classes in

TABLE II
TUNED PARAMETERS FOR THE GRADIENT BOOSTING CLASSIFIER

| Parameter | Value |
|-------------------|-------|
| learning_rate | 0.01 |
| max_depth | 9 |
| max_features | 11 |
| min_samples_leaf | 1 |
| min_samples_split | 7 |
| n_estimators | 1000 |
| subsample | 0.8 |

TABLE III
MEAN F1, F2 AND F0.5 SCORES FOR ALL FOUR CLASSIFIERS OVER 10 ROUNDS OF 5-FOLD CROSS VALIDATION.

| | RandomForest | ExtraTrees | GradientBoost | AdaBoost |
|------|--------------|------------|---------------|----------|
| F1 | 0.597687 | 0.588262 | 0.624193 | 0.624870 |
| F2 | 0.519876 | 0.539351 | 0.570319 | 0.574151 |
| F0.5 | 0.713685 | 0.656862 | 0.698756 | 0.691422 |

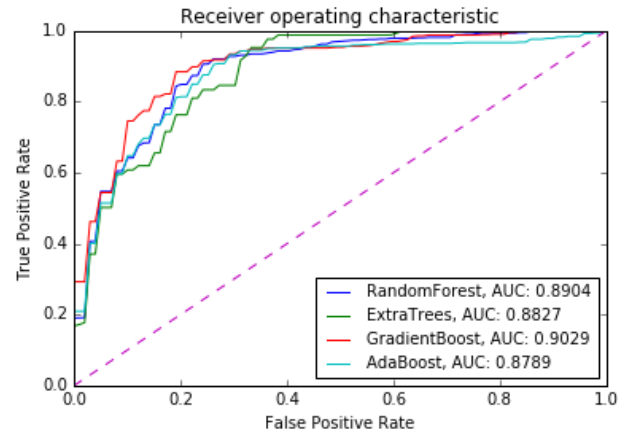


Fig. 2. Mean ROC curves and the mean AUC for all four classifiers over 10 rounds of 5-fold cross validation. All four classifiers perform similarly, with Extra Trees producing the best AUC. All four curves are slightly skewed towards the right of the plot, suggestive of poor recall. The F2 score (table III and figure 3 confirms this.

each fold. For all classifiers cross validation was repeated ten times, each with a new set of folds to ensure consistent results.

Figure 2 shows the mean ROC curve and mean AUC for each of the classifiers after cross validation. The performance of each classifier appears to be very similar. Notably the ROC curve for each type of classifier is shifted to the right of the graph, suggesting that they all exhibit a low recall rate.

Table III and figure 3 confirm this indication. Table III shows the F measure with a β parameter of 1, 2, and 0.5. Figure 3 shows a bar chart of the F2 scores in table III. The performance of all classifiers measured with the F2 score (which weights recall more highly than precision) is much lower in comparison to the F0.5 and F1 scores. This further confirms that all classifiers have a problem with recall.

B. Feature Engineering

In addition to the preprocessing steps outlined in II-A several combinations of new features were generated from the existing predictors. Firstly, as a large portion of the features

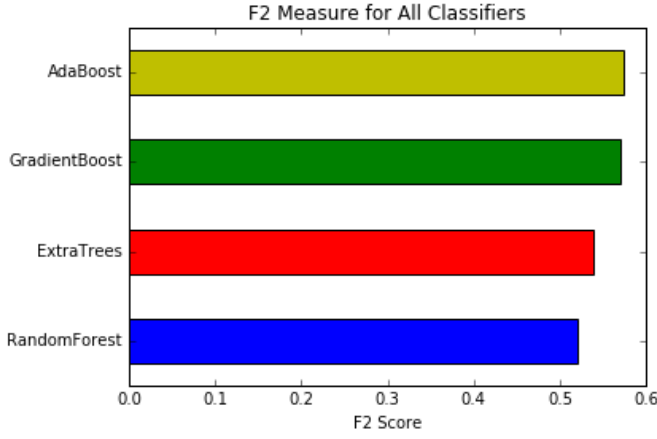


Fig. 3. Bar chart showing the F2 score for all classifiers taken from table III.

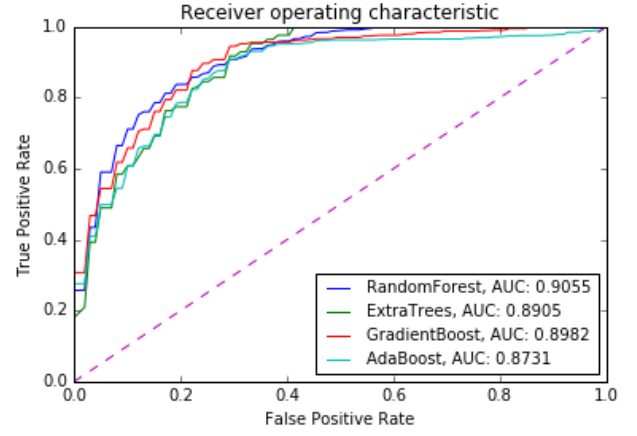


Fig. 5. ROC curves for each of the classifiers with the spirometry features.

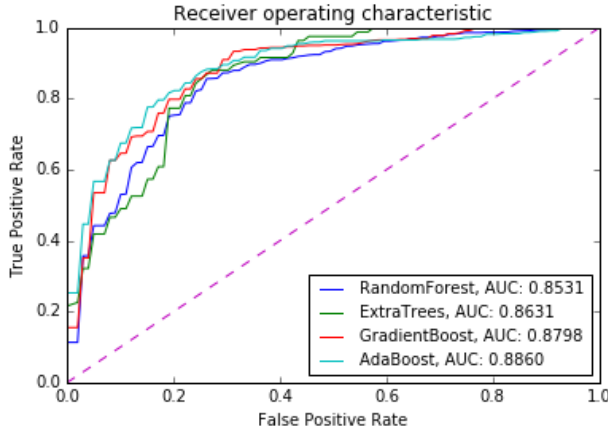


Fig. 4. ROC curves for each of the classifiers with the additional binary features. Performance is notably worse compared to the initial run.

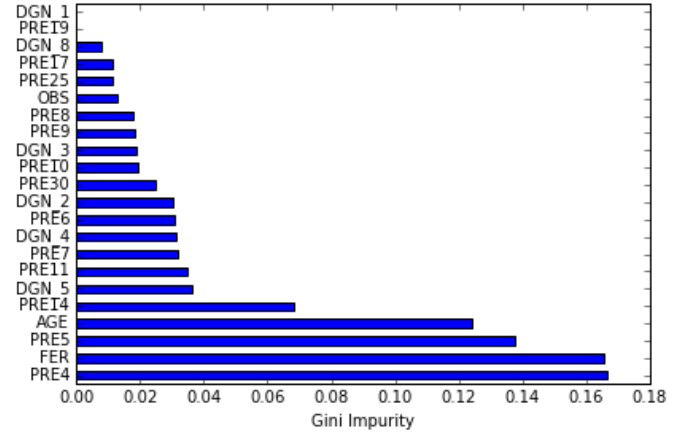


Fig. 6. Variable importance for each of the features with the spirometry features included.

are binary, a set of new features were created based on logical binary operators. The creation of the binary features is as follows: all pairs of binary features are enumerated. From each pair three new features are created by combining the pair using logical OR, AND and XOR.

Figure 4 shows the ROC curves for the same dataset but with the additional binary features appended. Table IV shows the corresponding F-scores for each classifier. From these results it is easy to see that all of the classifiers appear to perform worse with the new features. This may be due to their already limited contribution and that fact that the additional dimensionality is hindering progress.

The second modification to the original dataset is to create a couple of new features called FER and OBS. FER is the

FEV1/FVC ratio which is a spirometry measurement defined as $(FEV1/FVC) \cdot 100$ [9]. It is interpreted as the percentage of FVC expelled in the first second of a forced expiration. A ratio value below $<70\%$ can be suggestive of an obstructive disease. Using this information another feature (OBS) is generated from the ratio. OBS is a binary feature with value 1 when a patient has a ratio $<70\%$.

From figure 5 it can be seen that there is a slight improvement over the original ROC AUC scores using these additional spirometry features. This is backed up by plotting the feature importance obtained from training a random forest on the dataset (see 1). The two new features, particularly the FER feature are providing useful training information. This seems sensible as the new features are just combinations of existing well performing features.

TABLE IV
F SCORES FOR THE DATASET INCLUDING BINARY FEATURES

| | RandomForest | ExtraTrees | GradientBoost | AdaBoost |
|------|--------------|------------|---------------|----------|
| F1 | 0.532500 | 0.571757 | 0.611276 | 0.652767 |
| F2 | 0.459203 | 0.517610 | 0.543585 | 0.598086 |
| F0.5 | 0.644672 | 0.646104 | 0.707112 | 0.728810 |

TABLE V
F SCORES FOR THE DATASET INCLUDING SPIROMETRY FEATURES

| | RandomForest | ExtraTrees | GradientBoost | AdaBoost |
|------|--------------|------------|---------------|----------|
| F1 | 0.566668 | 0.607952 | 0.601360 | 0.619702 |
| F2 | 0.484577 | 0.562422 | 0.540406 | 0.560452 |
| F0.5 | 0.696171 | 0.670030 | 0.687827 | 0.699803 |

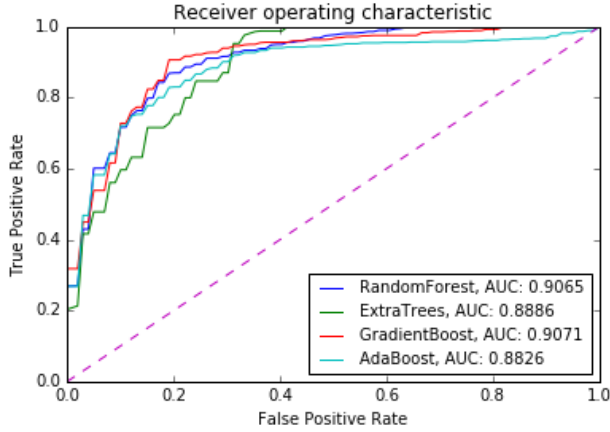


Fig. 7. ROC curves for each of the classifiers with the polynomial combination features.

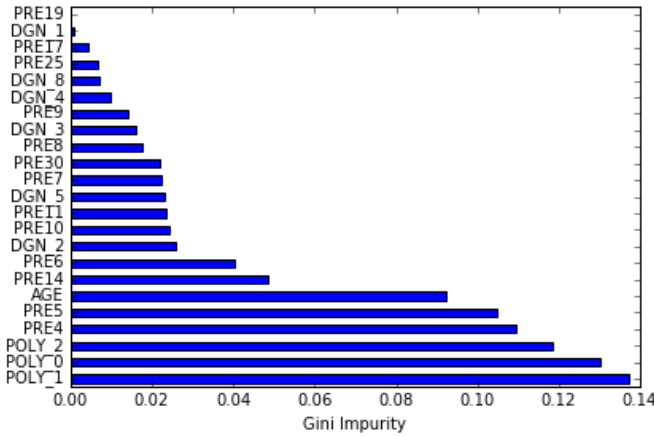


Fig. 8. Variable importance for each of the features with the polynomial combination features included.

Motivated by the results of the previous test, a selection of new features were created from all order 2 polynomial combinations of the two best predictors: PRE4 and PRE5. This means that the new features are of the form a^2 , ab , b^2 where a and b are PRE4 and PRE5 respectively.

Polynomial combinations led to the best result of the new features across all classifiers under cross validation. Figure 7 shows the ROC curves with the additional features included. The F-scores for each classifier are shown in table VI. The contribution of the new features can be seen in the variable importance plot (figure 8).

TABLE VI
F SCORES FOR THE DATASET INCLUDING POLYNOMIAL COMBINATION FEATURES

| | RandomForest | ExtraTrees | GradientBoost | AdaBoost |
|------|--------------|------------|---------------|----------|
| F1 | 0.596854 | 0.586917 | 0.617784 | 0.658825 |
| F2 | 0.508636 | 0.539122 | 0.563169 | 0.608835 |
| F0.5 | 0.733492 | 0.653059 | 0.692061 | 0.723445 |

TABLE VII
MEAN F1, F2 AND F0.5 SCORES FOR ALL CLASSIFIERS AFTER MONTE CARLO CROSS VALIDATION WITH SMOTE RESAMPLING WITH A RATIO OF 0.8

| | RandomForest | ExtraTrees | GradientBoost | AdaBoost |
|------|--------------|------------|---------------|----------|
| F1 | 0.608007 | 0.603050 | 0.628132 | 0.609344 |
| F2 | 0.636271 | 0.643457 | 0.652385 | 0.640834 |
| F0.5 | 0.585051 | 0.570976 | 0.608498 | 0.583632 |

C. Dataset Balancing

As mentioned in section II-A the thoracic surgery dataset is class imbalanced with only 28% of the dataset being of the positive class. One technique to combat class imbalance is to resample the dataset to put more emphasis on the known positive examples. A popular technique for resampling data is SMOTE [10]. SMOTE rebalances a dataset by creating new synthetic training to balance out the majority class using the nearest neighbours to a training instance. SMOTE is typically combined with under-sampling of the majority class to produce a final dataset that is re-weighted in favour of the minority class.

The results for the classifiers in part III-A shows that they have lower recall than precision. Rebalancing the dataset should show a decrease in precision and an increase in recall rate. This can be desirable in a dataset such as this where recall may be more important than precision. It is probably more desirable overestimate the number people who are likely to die from surgery than to achieve high precision.

SMOTE modified datasets cannot be validated using conventional k-fold cross validation. This is because the testing fold would contain synthetically generated training examples which are obviously not representative of the ground truth. Instead, in order to achieve a representative sample of performance, Monte Carlo cross-validation [11] is used. Before any resampling is applied, the data set is randomly split into a training and testing set. The split is stratified according to the class labels. All reported experiments use and 80/20 split. Resampling is then applied to the training dataset only, with the testing set remaining untouched. This process is then repeated for the desired number of iterations and the resulting performance measures are averaged. In all experiments the number of iterations performed was 50.

Figure 9 shows ROC curve and mean AUC scores for each of the classifiers using SMOTE with a resampling ratio of 0.8. Table VII shows the F-scores for each of the classifiers. Comparing this table to the results of III shows a clear difference in the F2 score. Recall weighted performance is now better both than F1 and F0.5. This improvement comes at the cost of a decrease in both the AUC and F0.5 measures. Increasing the oversampling ratio or under-sampling the majority class accentuates this effect.

IV. DISCUSSION

The experiments in section III have shown a variety of different approaches to predicting surgery survival with ensemble methods. Some of the best performance was achieved using the just the basic preprocessing steps outlined in section II-A.

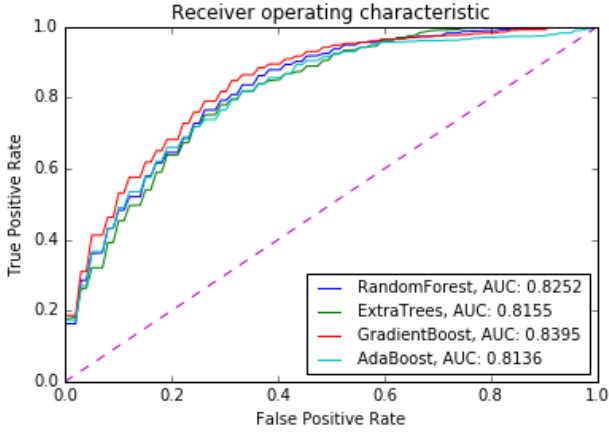


Fig. 9. ROC curves for all four classifiers with SMOTE oversampling with a ratio of 0.8. Each curve represents the average over 50 iterations of Monte Carlo validation. The ROC curves for all classifiers are less skewed compared to figure 2.

Looking at the initial performance evaluation (figure 2) it can be seen that all classifiers performed reasonable similarly with Gradient Boosting narrowly coming out ahead. The weakest performer was AdaBoost. Looking at the F-scores for each of the classifiers (table III) is more informative. All classifiers can be seen to perform comparable. Each performed weaker under the F2 measure which weights recall more highly than precision. It is the higher recall rate which is primarily driving the improvement of Gradient Boosting over the other classifiers in this trial.

Motivated by this baseline evaluation, this paper explored alternative feature representations through “feature engineering”. The first attempt was to create combinations of binary features from the existing dataset. This actually lead to worse results compared to the original preprocessing steps. None of the new binary features significantly contributed new information for the algorithm to work with. This probably meant that the increase in dimensionality out weighed any small gains delivered by the new representation. Each tree in the ensembles will only look at a limited number of features, so adding lots of redundant features is only likely to decrease performance while increasing training time. Note that an attempt was made to reduce the number of features by discarding the n weakest features using both variable importance and PCA, but neither method improved the results of this test.

The second feature engineering experiment was to derive a couple of features from spirometry theory. Both the FER and OBS features are directly derived from the existing predictors in the dataset. The bar chart in figure 6 shows that these features do appear to contribute some additional information. This is reflected in the AUC scores and the corresponding F-scores. The result for AUC is nearly identical between Gradient Boosting, Extra-Trees, and Random forests. The most successful (by a tiny margin) showed that Random forests performed the best in terms in AUC but this was only due to high precision. The F2 and F1 score are both much reduced in in the Random Forest trial. With these new features the best

candidate appears to be Extra-Trees which has an improvement across all three F-scores.

The third experiment involving feature engineering was to create a new batch of features by creating polynomial combinations of the features PRE4 and PRE5 which were shown to be the strongest predictors in figure 1. This led to some of the best AUC scores out of all the trails with two out of the four classifiers pushing into the 0.9 range. The variable importance plot (8) shows that the polynomial features are the most successful contributors. The F-scores show mixed results across the board. Notably AdaBoost in particular faired much better with polynomial features.

Finally an experiment was carried looking at improving the performance by resampling the dataset using SMOTE. While Random Forests and Extra-Trees already carry out bagging which already balances the dataset during training, this method could of potentially helped the performance of the boosting classifiers. Figure 9 shows a marked decrease in performance across all classifiers. This is probably due to the synthetic examples not realistically reflecting the distribution of positive examples in the dataset. What is more interesting is the fact that the F2 scores for each classifier are improved by applying SMOTE but the precision is dramatically hindered. This is could be due to the synthetic examples “expanding” the region around positive examples which the algorithm considers to be positive. This is probably not representative of the true decision boundary, but has the effect of increasing recall as more examples are likely to land with the expanded positive region. While this test resulted in much worse performance it could still be of interest. In predicting thoracic surgery survival it is more desirable to have high recall than high precision.

V. CONCLUSIONS

In conclusion this paper examined the effect of four different ensemble methods on a variety of engineered features. The effect of resampling the dataset was also explored. The final classifier used on the unseen testing data for submission was trained using both the additional polynomial and spirometry features. This lead to best performance with the Gradient Boosting classifier. This classifier had a final AUC after cross validation of 0.903, just under the score for Random Forests, but it had a much more balanced set of F-scores compared to Random Forests.

While this is one of the best AUC scores achieved across all experiments there is clearly some room from improvement. One area of improvement worth exploring would be to look at automated methods of feature selection. One such method could be recursive feature elimination in conjunction with a model that estimates feature relevance (such as random forests). Alternatively a non-linear dimensionality reduction technique could be used to find projections of the feature space onto a lower dimensional embedding. This could be particularly beneficial in the case of the binary features where the feature space is relatively much larger.

Another avenue for exploration would be to look at a different branch of algorithms. For example a penalised SVM could be experimented with. The original authors of [1]

proposed a boosted SVM which performs reasonable well on the expanded dataset. Any alternative classifier will probably benefit from some form of bagging or resampling more than the ensemble methods.

The experiments in this paper show that any further predictive progress appears to be hindered by the low recall rate. This is a common trade-off in machine learning. Implementing this system in the real world would most likely require favouring a lower F0.5 score for a higher F2 for safety reasons. Any progress beyond the AUC achieved in this paper is likely to require a combination of further creative feature engineering and a good mix of bagging/resampling.

REFERENCES

- [1] M. Zięba, J. M. Tomczak, M. Lubicz, and J. Świątek, “Boosted svm for extracting rules from imbalanced data in application to prediction of the post-operative life expectancy in the lung cancer patients,” *Applied soft computing*, vol. 14, pp. 99–108, 2014.
- [2] L. Breiman, “Random forests,” *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [3] P. Geurts, D. Ernst, and L. Wehenkel, “Extremely randomized trees,” *Machine learning*, vol. 63, no. 1, pp. 3–42, 2006.
- [4] Y. Freund and R. E. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of computer and system sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [5] A. Natekin and A. Knoll, “Gradient boosting machines, a tutorial,” *Frontiers in neurorobotics*, vol. 7, 2013.
- [6] “Thoracic Surgery Data Set,” <http://archive.ics.uci.edu/ml/datasets/Thoracic+Surgery+Data>, accessed: 2016-04-29.
- [7] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [8] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, “Scikit-learn: Machine learning in python,” *The Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [9] “Spirometry, Patient Website,” <http://patient.info/doctor/spirometry-pro>, accessed: 2016-05-04.
- [10] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, “Smote: synthetic minority over-sampling technique,” *Journal of artificial intelligence research*, pp. 321–357, 2002.
- [11] W. Dubitzky, M. Granzow, and D. P. Berrar, *Fundamentals of data mining in genomics and proteomics*. Springer Science & Business Media, 2007.

APPENDIX A THIRD PARTY LIBRARIES

The code use to produce the results in the paper rely upon a number of different third party libraries. The libraries used and the relevant version of each is show in the table below. lease note that the *UnbalancedData* library is not currently directly available through *pip* by default and so must be installed directly from the GitHub repository using the following command:

```
1 pip install git+https://github.com/fmfn/UnbalancedDataset
```

TABLE VIII
THIRD PARTY LIBRARIES AND THE ACCOMPANYING VERSIONS USED IN ALL THE FOLLOWING CODE SAMPLES.

| Name | Version |
|-------------------|---------|
| pandas | 0.18.0 |
| sklearn | 0.17.1 |
| UnbalancedDataset | 0.1 |
| matplotlib | 1.5.1 |
| numpy | 1.11.0 |
| scipy | 0.17.0 |

APPENDIX B PYTHON SCRIPT FOR FINAL CLASSIFIER

This appendix contains the python script for creating the final classifier used to make the predictions on the test dataset for this assignment. This will save a CSV file in the data folder called *submission.csv*.

```
1 import pandas as pd
2 import numpy as np
3 from sklearn import preprocessing
4 from sklearn.pipeline import Pipeline
5 from sklearn.ensemble import GradientBoostingClassifier
6
7 df = pd.DataFrame.from_csv("../data/train_risk.csv", index_col=False)
8 test = pd.DataFrame.from_csv("../data/test_risk.csv", index_col=False)
9 X, y = df[df.columns[:-1]], df[df.columns[-1]]
10
11
12 def encode_onehot(x_data, column_name, digitize=False):
13     """ Encode a catagorical column from a data frame into a
14     data frame of one hot features
15     """
16     data = x_data[[column_name]]
17
18     if digitize:
19         data = np.digitize(data, np.arange(data.min(), data.max(), 10))
20
21     enc = preprocessing.OneHotEncoder()
22     features = enc.fit_transform(data).toarray()
23     names = ['%s_%d' % (column_name, i) for i in enc.active_features_]
24     features = pd.DataFrame(features, columns=names, index=x_data.index)
25     return features
26
27
28 def create_spiro_features(x_data):
29     """ Create spriometry based features """
30     # create new feature FER
31     # this is the raito of FEV1 and FVC
32     FER = (x_data.PRE5 / x_data.PRE4) * 100
33     FER.index = x_data.index
34
35     # create a new feature OBS
36     # this is whether the instance has a FER below 70%
37     # which implies an obstructive disease.
38     OBS = pd.Series(np.zeros(x_data.AGE.shape))
39     OBS.index = x_data.index
40     OBS.loc[FER < 70] = 1.0
41
42     spiro = pd.concat([FER, OBS], axis=1)
43     spiro.columns = ['FER', 'OBS']
44     return spiro
45
46
```



```

47 def create_poly_features(x_data, names):
48     """ Create new features base on Polynomials of the original best two predictors """
49     poly = preprocessing.PolynomialFeatures(2, include_bias=False, interaction_only=True)
50     poly_features = pd.DataFrame(poly.fit_transform(x_data[names]), index=x_data.index)
51     poly_features.columns = ["POLY_%d" % i for i in poly_features.columns]
52     return poly_features
53
54
55 def preprocess(x_data, y_data=None):
56     # drop zero var PRE32
57     Xp = x_data.drop("PRE32", axis=1)
58
59     # remove outliers
60     if y_data is not None:
61         mask = Xp.PRE5 < 30
62         Xp = Xp.loc[mask]
63         Yp = y_data.copy()
64         Yp = Yp.loc[mask]
65     else:
66         Yp = None
67
68     # encode catagorical data as one hot vectors
69     one_hot_names = ["DGN"]
70     encoded = map(lambda name: encode_onehot(Xp, name), one_hot_names)
71     # combine into a single data frame
72     new_features = pd.concat(encoded, axis=1)
73
74     # drop the catagorical variables that have been encoded
75     Xp.drop(["DGN"], inplace=True, axis=1)
76     # add new features
77     Xp = pd.concat([Xp, new_features], axis=1)
78
79     return Xp, Yp
80
81 Xp, Yp = preprocess(X, y)
82
83 scaler = preprocessing.StandardScaler()
84
85 gbc_params = {
86     'min_samples_leaf': 1,
87     'min_samples_split': 7,
88     'max_depth': 9,
89     'max_features': 11,
90     'subsample': 0.8,
91     'n_estimators': 1000,
92     'learning_rate': 0.01
93 }
94
95 gbc = GradientBoostingClassifier(**gbc_params)
96 gbc_pipe = Pipeline([('scaler', scaler), ('GradientBoostingClassifier', gbc)])
97
98 model = {'name': 'GradientBoosting', 'model': gbc_pipe}
99
100 # Create training features
101 spiro_features = create_spiro_features(Xp)
102 poly_features = create_poly_features(Xp, ['PRE4', 'PRE5'])
103 Xp_all = pd.concat([Xp, poly_features, spiro_features], axis=1)
104 Xp_all.drop(['DGN_1', 'DGN_8'], axis=1, inplace=True)
105
106 # Create testing features
107 Xtest, _ = preprocess(test, y_data=None)
108 Xtest = Xtest.drop('test_id', axis=1)
109
110 test_spiro_features = create_spiro_features(Xtest)
111 test_poly_features = create_poly_features(Xtest, ['PRE4', 'PRE5'])
112 Xtest = pd.concat([Xtest, test_spiro_features, test_poly_features], axis=1)
113
114 # Build model
115 final_model = model['model']
116 final_model.fit(Xp_all, Yp)
117 predicted_prob = pd.Series(final_model.predict_proba(Xtest)[: , 1])
118 predicted_label = pd.Series(final_model.predict(Xtest))
119
120 final_submission = pd.concat([test.test_id, predicted_label, predicted_prob], axis=1)
121 final_submission.columns = ['test_id', 'predicted_label', 'predicted_output']
122 final_submission.to_csv('../data/submission.csv', index=False)

```

../src/classify.py

APPENDIX C

IPYTHON NOTEBOOK AND ADDITIONAL PYTHON MODULES FOR ANALYSIS, TRAINING AND TUNING

This listing shows the contents of the analysis IPython notebook as a python script. For a better formatted version of this code install IPython and open the Analysis.ipynb file provided with the assignment submission. This IPython notebook contains all of the code for analysing the data, tuning the algorithms, and performing both stratified k-fold and Mote Carlo cross validation. Two additional modules (*pipeline* and *roc_analysis*) are also provided as listing as the end of this section.

```

1
2 # coding: utf-8
3
4 # In[1]:
5
6 get_ipython().magic(u'matplotlib inline')
7 import pandas as pd
8 import numpy as np
9 import matplotlib.pyplot as plt
10 from sklearn.pipeline import Pipeline
11 from sklearn.svm import SVC
12 from sklearn.ensemble import ExtraTreesClassifier, AdaBoostClassifier, GradientBoostingClassifier,
   RandomForestClassifier
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.tree import DecisionTreeClassifier
15 import pipeline
16 from roc_analysis import ROCAnalysisScorer
17
18
19 # ### Loading the Datasets
20
21 # In[197]:
22
23 df = pd.DataFrame.from_csv("../data/train_risk.csv", index_col=False)
24 test = pd.DataFrame.from_csv("../data/test_risk.csv", index_col=False)
25 X, y = df[df.columns[:-1]], df[df.columns[-1]]
26
27
28 # ## Analysing the Data
29 #
30 # Looking at the difference between the number of positive and negative samples in the dataset shows that
   there are more negative examples than positive examples. Only 28% of all samples are of the positive
   class.
31
32 # In[3]:
33
34 def class_balance_summary(y):
35     """ Summarise the imbalance in the dataset """
36     total_size = y.size
37     negative_class = y[y == 0].size
38     positive_class = y[y > 0].size
39     ratio = positive_class / float(positive_class + negative_class)
40
41     print "Total number of samples: %d" % total_size
42     print "Number of positive samples: %d" % positive_class
43     print "Number of negative samples: %d" % negative_class
44     print "Ratio of positive to total number of samples: %.2f" % ratio
45
46
47 class_balance_summary(y)
48
49
50 # Some initial observations about the data before it is preprocessed:
51 # - PRE32 is all zeros. This can be removed
52 # - PRE14 looks catagorical. Should be split into multiple binary variables
53 # - DGN looks catagorical. As above.
54 # - PRE5 looks to have some outliers. See box plot below. Potentially remove or split into two extra
   variable?
55
56 # In[4]:
57
58 X.head()
59

```

```

60
61 # Box plot below shows the outliers in PRE5. It is worth noting that all of these outliers are of the
    negative class. This variable is the volume that can be exhaled in one second given full inhalation. It
    is likely that these values are therefore errors in reporting as it is unlikely that humans can exhale
    such a large volume so quickly.
62
63 # In[5]:
64
65 # X.PRE5.plot(kind='box')
66 X.PRE5.plot(kind='box')
67 print y[X.PRE5 > 30 ]
68
69
70 # ## Preprocessing
71 #
72 # Create a new matrix of preprocessed features. This will encode catagorical data as one hot vectors,
    remove outliers, and normalise the data.
73
74 # In[198]:
75
76 from sklearn import preprocessing
77
78 def encode_onehot(x_data, column_name, digitize=False):
79     """ Encode a catagorical column from a data frame into a data frame of one hot features"""
80     data = x_data[[column_name]]
81
82     if digitize:
83         data = np.digitize(data, np.arange(data.min(), data.max(), 10))
84
85     enc = preprocessing.OneHotEncoder()
86     features = enc.fit_transform(data).toarray()
87     names = ['%s_%d' % (column_name, i) for i in enc.active_features_]
88     features = pd.DataFrame(features, columns=names, index=x_data.index)
89     return features
90
91
92 def preprocess(x_data, y_data=None):
93     # drop zero var PRE32
94     Xp = x_data.drop("PRE32", axis=1)
95
96     # remove outliers
97     if y_data is not None:
98         mask = Xp.PRE5 < 30
99         Xp = Xp.loc[mask]
100         Yp = y_data.copy()
101         Yp = Yp.loc[mask]
102     else:
103         Yp = None
104
105     # encode catagorical data as one hot vectors
106     one_hot_names = ["DGN"]
107     encoded = map(lambda name: encode_onehot(Xp, name), one_hot_names)
108     #combine into a single data frame
109     new_features = pd.concat(encoded, axis=1)
110
111     # drop the catagorical variables that have been encoded
112     Xp.drop(["DGN"], inplace=True, axis=1)
113     # add new features
114     Xp = pd.concat([Xp, new_features], axis=1)
115
116     return Xp, Yp
117
118 Xp, Yp = preprocess(X, y)
119 Xp.head()
120
121
122 # Measure the effectiveness of each feature using the variable importance measure from a Random Forest
123
124 # In[10]:
125
126 from sklearn.ensemble import RandomForestClassifier
127 from sklearn.preprocessing import StandardScaler
128
129 def measure_importance(x_data, y_data):
130     rf_selector = RandomForestClassifier(criterion='gini', class_weight='balanced')
131     rf_selector.fit(StandardScaler().fit_transform(x_data), y_data)
132     feature_importance = pd.Series(rf_selector.feature_importances_, index=x_data.columns).sort_values(

```

```

    ascending=False)
    feature_importance.plot(kind='bar')
    return feature_importance
133
134
135
136 feature_importance = measure_importance(Xp, Yp)
137 Xp.drop(feature_importance[feature_importance == 0].index, inplace=True, axis=1)
138
139
140 # In[11]:
141
142 feature_importance.plot(kind='barh')
143 plt.xlabel('Gini Impurity')
144 plt.tight_layout()
145 plt.savefig("img/feature_importance.png")
146
147
148 # The numerical features appear to be the most important ones. Plot a scatter plot matrix to see how the
    how the correlate with each other
149
150 # In[12]:
151
152 pd.tools.plotting.scatter_matrix(Xp[['PRE4', 'PRE5', 'AGE']], c=Yp)
153
154
155 # ## Tuning Model Parameters
156 #
157 # Given the current status of the data tune the model parameters to it before we evaluate the overall
    performance. Note that all of the tuning presented here is orientated towards obtaining the highest AUC
    score. Other metrics might be more desirable given the problem domain, but AUC is the measurement used
    for assignment points.
158
159 # In[13]:
160
161 from sklearn import cross_validation
162 skf = cross_validation.StratifiedKFold(Yp, n_folds=5)
163
164
165 # ### Random Forest Tuning
166 # Run a grid search over a range of parameters for a Random Forest. The dataset is small enough that we can
    do them all at once. '''n_estimators''' is neglected because this should always improve as it is
    increased so we should attempt to make it as large as possible subject to lack of improvement
167
168 # In[3835]:
169
170 param_grid = {"max_depth": range(2, 20, 3),
171               "max_features": range(2, 20, 3),
172               "min_samples_split": range(1, 5),
173               "min_samples_leaf": range(1, 5),
174               }
175
176 rf = RandomForestClassifier(class_weight='balanced', n_estimators=50, random_state=50)
177 rf_clf = grid_search.GridSearchCV(rf, param_grid, n_jobs=-1, cv=skf, scoring='roc_auc')
178 rf_clf.fit(Xp, Yp)
179
180
181 # In[3836]:
182
183 print rf_clf.best_params_
184
185
186 # Now take a look at the number of estimators and see where performance begins to level off.
187
188 # In[3838]:
189
190 param_grid = {"n_estimators": range(50, 500, 50)}
191 const_params = {'max_features': 1, 'min_samples_split': 1, 'max_depth': 16, 'min_samples_leaf': 1}
192
193 rf = RandomForestClassifier(class_weight='balanced', n_estimators=50, random_state=50, **const_params)
194 rf_clf2 = grid_search.GridSearchCV(rf, param_grid, n_jobs=-1, cv=skf, scoring='roc_auc')
195 rf_clf2.fit(Xp, Yp)
196
197
198 # The best parameters for '''n_estimators''' levels off after around 300 estimators
199
200 # In[3840]:
201
202 plt.plot([d[0]['n_estimators'] for d in rf_clf2.grid_scores_], [d[1] for d in rf_clf2.grid_scores_])

```

```

203 print rf_clf2.best_params_
204 print rf_clf2.best_score_
205
206
207 # In[3725]:
208
209 rf_clf2.best_estimator_.get_params()
210
211
212 # ### Gradient Boosting Tuning
213 #
214 # Gradient boosting is difficult to tune effectively. [This guide](http://www.analyticsvidhya.com/blog
    /2016/02/complete-guide-parameter-tuning-gradient-boosting-gbm-python/) suggests starting by fixing the
    learning rate and number of estimators to a relatively low number in order to tune the other
    hyperparameters. After they are optimised the learning rate is gradually lowered and the number of
    estimators increased until we find convergence on the optimum parameters
215
216 # In[3587]:
217
218 param_grid = [
219     {'n_estimators': range(20,150,10)}
220 ]
221
222 const_params = {'learning_rate': 0.1, 'min_samples_split': 1, 'min_samples_leaf': 3, 'max_depth': 8, '
    max_features': 'sqrt', 'subsample': 0.8}
223 gbc = GradientBoostingClassifier(random_state=50, **const_params)
224
225 gbc_clf = grid_search.GridSearchCV(gbc, param_grid, cv=skf, scoring='roc_auc')
226 gbc_clf.fit(Xp, Yp)
227
228
229 # '''n_estimators''' plateaus at around 100, so we'll use this instead of the optimum as less trees ==
    quicker training and we'll need to decrease the learning rate and increase the number of trees later in
    the tuning anyway.
230
231 # In[3588]:
232
233 plt.plot([d[0]['n_estimators'] for d in gbc_clf.grid_scores_], [d[1] for d in gbc_clf.grid_scores_])
234 print gbc_clf.best_params_
235
236
237 # Now tune the '''max_depth''' and the '''min_samples_split''' parameters.
238
239 # In[3594]:
240
241 const_params = {'n_estimators':100,
242                 'learning_rate': 0.1,
243                 'min_samples_leaf': 3,
244                 'max_features': 'sqrt',
245                 'subsample': 0.8
246                 }
247
248 param_grid = [
249     {'max_depth': range(5,16,2), 'min_samples_split': range(1, 20, 3)}
250 ]
251
252
253 gbc = GradientBoostingClassifier(random_state=50, **const_params)
254 gbc_clf = grid_search.GridSearchCV(gbc, param_grid, cv=skf, scoring='roc_auc')
255 gbc_clf.fit(Xp, Yp)
256
257
258 # In[3595]:
259
260 print gbc_clf.best_params_
261 gbc_clf.grid_scores_
262
263
264 # Now train '''max_features'''
265
266 # In[3600]:
267
268 const_params = {'n_estimators':100,
269                 'learning_rate': 0.1,
270                 'min_samples_leaf': 3,
271                 'max_features': 'sqrt',
272                 'max_depth': 9,

```

```

273         'min_samples_split': 7,
274         'subsample': 0.8
275     }
276
277 param_grid = [
278     {'max_features': range(5, 20, 2)}
279 ]
280
281 gbc = GradientBoostingClassifier(random_state=50, **const_params)
282 gbc_clf = grid_search.GridSearchCV(gbc, param_grid, cv=skf, scoring='roc_auc')
283 gbc_clf.fit(Xp, Yp)
284
285
286 # In[3601]:
287
288 plt.plot([d[0]['max_features'] for d in gbc_clf.grid_scores_], [d[1] for d in gbc_clf.grid_scores_])
289 print gbc_clf.best_params_
290
291
292 # Now train to tune the “subsample” rate.
293
294 # In[3603]:
295
296 const_params = {'n_estimators': 100,
297                 'learning_rate': 0.1,
298                 'min_samples_leaf': 3,
299                 'max_features': 'sqrt',
300                 'max_depth': 9,
301                 'min_samples_split': 7,
302                 'max_features': 11,
303                 'subsample': 0.8
304                 }
305
306
307 param_grid = [
308     {'subsample': [0.6, 0.7, 0.75, 0.8, 0.85, 0.9]}
309 ]
310
311
312 gbc = GradientBoostingClassifier(random_state=50, **const_params)
313 gbc_clf = grid_search.GridSearchCV(gbc, param_grid, cv=skf, scoring='roc_auc')
314 gbc_clf.fit(Xp, Yp)
315
316
317 # In[3604]:
318
319 plt.plot([d[0]['subsample'] for d in gbc_clf.grid_scores_], [d[1] for d in gbc_clf.grid_scores_])
320 print gbc_clf.best_params_
321
322
323 # Now cross validate with all the parameters set:
324
325 # In[3614]:
326
327 const_params = {
328     'min_samples_leaf': 1,
329     'min_samples_split': 7,
330     'max_depth': 9,
331     'max_features': 11,
332     'subsample': 0.8
333 }
334
335 param_grid = [
336     {'n_estimators': [100], 'learning_rate': [0.1]},
337     {'n_estimators': [200], 'learning_rate': [0.05]},
338     {'n_estimators': [1000], 'learning_rate': [0.01]},
339     {'n_estimators': [1500], 'learning_rate': [0.005]},
340 ]
341
342 gbc = GradientBoostingClassifier(random_state=50, **const_params)
343
344 gbc_clf = grid_search.GridSearchCV(gbc, param_grid, cv=skf, scoring='roc_auc')
345 gbc_clf.fit(Xp, Yp)
346
347
348 # In[3718]:
349

```



```

350 print gbc_clf.best_params_
351 gbc_clf.grid_scores_
352
353
354 # In[3854]:
355
356 p = pd.DataFrame(gbc_clf.best_estimator_.get_params(), index=['Value']).T
357 p.index.name = "Parameter"
358 print p.to_latex()
359
360
361 # ### AdaBoost Tuning
362 # Perhaps the easiest train due to a fairly limited number of parameters. Adjusting the ‘‘max_depth’’
363 # suggests that 4 appears to be roughly the best option for the depth of the decision trees.
364
365 # In[3764]:
366
367 param_grid = {"n_estimators": range(50, 1000, 50), 'learning_rate': [0.1, 0.5, 0.01, 0.005]}
368
369 dt = DecisionTreeClassifier(class_weight='balanced', max_depth=4)
370 adb = AdaBoostClassifier(dt)
371 adb_clf = grid_search.GridSearchCV(adb, param_grid, n_jobs=-1, cv=skf, scoring='roc_auc')
372 adb_clf.fit(Xp, Yp)
373
374 # In[3781]:
375
376 print adb_clf.best_params_
377 print adb_clf.best_score_
378 adb_clf.grid_scores_
379
380
381 # ### Extremely Random Trees Tuning
382 #
383 # This is very similar to Random Forests. In fact we will start with the same parameter set for the grid
384 # search.
385
386 # In[3800]:
387
388 param_grid = {"max_depth": range(2, 20, 3),
389               "max_features": range(2, 20, 3),
390               "min_samples_split": range(1, 5),
391               "min_samples_leaf": range(1, 5),
392               }
393 etc = ExtraTreesClassifier(class_weight='balanced', bootstrap=True, n_estimators=50, random_state=50)
394 etc_clf = grid_search.GridSearchCV(etc, param_grid, n_jobs=-1, cv=skf, scoring='roc_auc')
395 etc_clf.fit(Xp, Yp)
396
397 # In[3801]:
398
399 print etc_clf.best_params_
400 print etc_clf.best_score_
401
402
403 # Now check increasing the number of estimators and find the drop off point
404
405 # In[3805]:
406
407 param_grid = {"n_estimators": range(50, 500, 50)}
408 const_params = {'max_features': 16, 'min_samples_split': 1, 'max_depth': 19, 'min_samples_leaf': 1}
409
410 etc = ExtraTreesClassifier(class_weight='balanced', bootstrap=True, random_state=50, **const_params)
411 etc_clf2 = grid_search.GridSearchCV(etc, param_grid, n_jobs=-1, cv=skf, scoring='roc_auc')
412 etc_clf2.fit(Xp, Yp)
413
414
415 # In[3806]:
416
417 plt.plot([d[0]['n_estimators'] for d in etc_clf2.grid_scores_], [d[1] for d in etc_clf2.grid_scores_])
418 print etc_clf2.best_params_
419 print etc_clf2.best_score_
420
421
422 # In[3807]:
423
424 etc_clf2.best_estimator_.get_params()

```

```

425
426
427 # ## Model Performance
428 # Test the performance of each of the models on the preprocessed dataset before trying any more complicated
    feature engineering/resampling. This should give us some rough baseline AUC measures to work with.
    Firstly, set up the models. This creates a set of pipelines for each of the models we want to use.
429
430 # In[178]:
431
432 scaler = preprocessing.StandardScaler()
433
434 # set up classifier objects
435 knn = KNeighborsClassifier(n_neighbors=5, weights='distance')
436 dct = DecisionTreeClassifier(class_weight='balanced', max_depth=4)
437 abt = AdaBoostClassifier(dct, n_estimators=400, learning_rate=0.5)
438
439 gbc_params = {
440     'min_samples_leaf': 1,
441     'min_samples_split': 7,
442     'max_depth': 9,
443     'max_features': 11,
444     'subsample': 0.8,
445     'n_estimators': 1000,
446     'learning_rate': 0.01
447 }
448 gbc = GradientBoostingClassifier(**gbc_params)
449
450 exf_params = {
451     'bootstrap': False,
452     'class_weight': 'balanced',
453     'criterion': 'gini',
454     'max_depth': 19,
455     'max_features': 16,
456     'max_leaf_nodes': None,
457     'min_samples_leaf': 1,
458     'min_samples_split': 1,
459     'min_weight_fraction_leaf': 0.0,
460     'n_estimators': 200,
461     'n_jobs': 1,
462     'oob_score': False,
463     'random_state': 50,
464     'verbose': 0,
465     'warm_start': False
466 }
467
468
469 exf = ExtraTreesClassifier(**exf_params)
470
471 rf_params = {
472     'bootstrap': True,
473     'class_weight': 'balanced',
474     'criterion': 'gini',
475     'max_depth': 16,
476     'max_features': 1,
477     'max_leaf_nodes': None,
478     'min_samples_leaf': 1,
479     'min_samples_split': 1,
480     'min_weight_fraction_leaf': 0.0,
481     'n_estimators': 300
482 }
483 rf_balanced = RandomForestClassifier(**rf_params)
484
485 # create pipelines for each model
486 abt_pipe = Pipeline([('scaler', scaler), ('AdaBoost', abt)])
487 exf_pipe = Pipeline([('scaler', scaler), ('ExtraTrees', exf)])
488 gbc_pipe = Pipeline([('scaler', scaler), ('GradientBoostingClassifier', gbc)])
489 rfs_pipe = Pipeline([('scaler', scaler), ('RandomForest', rf_balanced)])
490
491 # create list of model data
492 models = [
493     {'name': 'AdaBoost', 'model': abt_pipe},
494     {'name': 'ExtraTrees', 'model': exf_pipe},
495     {'name': 'RandomForest', 'model': rfs_pipe},
496     {'name': 'GradientBoost', 'model': gbc_pipe},
497 ]
498
499 # set the same training set for all models.

```

```

500 # this is just the preprocessed dataset.
501 for model in models:
502     model['train_data'] = (Xp, Yp)
503
504
505 # Define some useful helper functions for summarising the results of k-fold/monte carlo cross validation
506
507 # In[150]:
508
509 def f_score_summary(scorers):
510     """ Create a summary of the average f-scores for all folds/trials """
511     series = []
512     columns = []
513     for key, scorer in scorers.iteritems():
514         f_scores = [np.mean(scorer.f1scores_), np.mean(scorer.f2scores_), np.mean(scorer.fhalf_scores_)]
515         s = pd.Series(f_scores, index=['F1', 'F2', 'F0.5'])
516         series.append(s)
517         columns.append(key)
518
519     frame = pd.concat(series, axis = 1)
520     frame.columns = columns
521     return frame
522
523 def summarise_scorers(scorers):
524     """ Create a summary of the scorers AUCs for all folds/trials """
525     names = [name for name in scorers.keys()]
526     aucs = [scorer.aucs_ for scorer in scorers.values()]
527     aucs = pd.DataFrame(np.array(aucs).T, columns=names)
528     return aucs.describe()
529
530
531 # Perform n iterations of k fold cross validation. Here I am using 10 iterations and 5 folds at each
532     iteration.
533
534 # In[151]:
535
536 scorers = pipeline.repeated_cross_fold_validation(models, n=10, k=5)
537
538 # Plot an ROC curve and the mean AUCs.
539
540 # In[153]:
541
542 for key, scorer in scorers.iteritems():
543     scorer.plot_roc_curve(mean_label=key, mean_line=True, show_all=False)
544
545 plt.plot(np.arange(0,1,1, 0.1), np.arange(0,1,1, 0.1), '—')
546 plt.savefig("img/roc_cv.png")
547
548
549 # Plot bar chart of the F2 scores
550
551 # In[154]:
552
553 f_scores = f_score_summary(scorers)
554 ax = f_scores.loc['F2'].plot(kind='barh', title='F2 Measure for All Classifiers', color=['b', 'r', 'g', 'y'])
555 ax.set_xlabel('F2 Score')
556 plt.tight_layout()
557 plt.savefig('img/f2_score.png')
558
559
560 # Summarise the F scores
561
562 # In[155]:
563
564 f_scores = f_score_summary(scorers)
565 print f_scores.to_latex()
566 f_scores
567
568
569 # ## Feature Engineering
570 #
571 # Test creating some new features based on combinations of existing ones in the dataset. Cross validate
572     each set of new features to see if it improves performance.
573
574 # ### Binary Features

```

```

574
575 # In[156]:
576
577 import itertools
578
579 def binary_combinations(x_data, names):
580     name_pairs = itertools.combinations(names, 2)
581     features = []
582     for a_name, b_name in name_pairs:
583         a, b = x_data[a_name], x_data[b_name]
584         features.append(np.logical_xor(a, b).astype(int))
585         features.append(np.logical_and(a, b).astype(int))
586         features.append(np.logical_or(a, b).astype(int))
587
588     return pd.DataFrame(np.array(features).T, index=x_data.index)
589
590 binary_features = binary_combinations(Xp, ['PRE7', 'PRE8', 'PRE9', 'PRE10', 'PRE11', 'PRE17', 'PRE30'])
591 Xp_binary = pd.concat([Xp, binary_features], axis=1)
592 feature_importance = measure_importance(Xp_binary, Yp)
593 Xp_binary.drop(feature_importance[feature_importance == 0].index, inplace=True, axis=1)
594
595
596 # In[157]:
597
598 for model in models:
599     model['train_data'] = (Xp_binary, Yp)
600
601 scorers = pipeline.repeated_cross_fold_validation(models, n=10, k=5)
602
603
604 # In[158]:
605
606 get_ipython().magic(u'matplotlib inline')
607 for key, scorer in scorers.iteritems():
608     scorer.plot_roc_curve(mean_label=key, mean_line=True, show_all=False)
609
610 plt.plot(np.arange(0,1.1, 0.1), np.arange(0,1.1, 0.1), '—')
611 plt.savefig("img/roc_binary_features.png")
612
613
614 # In[159]:
615
616 f_scores = f_score_summary(scorers)
617 print f_scores.to_latex()
618 f_scores
619
620
621 # ### Spirometry Based Features
622
623 # In[160]:
624
625 def create_spiro_features(x_data):
626     # create new feature FER
627     # this is the ratio of FEV1 and FVC
628     FER = (x_data.PRE5 / x_data.PRE4) * 100
629     FER.index = x_data.index
630
631     # create a new feature OBS
632     # this is whether the instance has a FER below 70%
633     # which implies an obstructive disease.
634     OBS = pd.Series(np.zeros(x_data.AGE.shape))
635     OBS.index = x_data.index
636     OBS.loc[FER < 70] = 1.0
637
638     spiro = pd.concat([FER, OBS], axis=1)
639     spiro.columns = ['FER', 'OBS']
640     return spiro
641
642
643 # In[161]:
644
645 spiro_features = create_spiro_features(Xp)
646 Xp_spiro = pd.concat([Xp, spiro_features], axis=1)
647 feature_importance = measure_importance(Xp_spiro, Yp)
648 Xp_spiro.drop(feature_importance[feature_importance == 0].index, inplace=True, axis=1)
649
650

```

```

651 # In[162]:
652
653 for model in models:
654     model['train_data'] = (Xp_spiro, Yp)
655
656 scorers = pipeline.repeated_cross_fold_validation(models, n=10, k=5)
657
658 # In[163]:
659
660 get_ipython().magic(u'matplotlib inline')
661 for key, scorer in scorers.iteritems():
662     scorer.plot_roc_curve(mean_label=key, mean_line=True, show_all=False)
663
664 plt.plot(np.arange(0,1.1, 0.1), np.arange(0,1.1, 0.1), '—')
665 plt.savefig("img/roc_spiro_features.png")
666
667 # In[164]:
668
669 feature_importance.plot(kind='barh')
670 plt.xlabel('Gini Impurity')
671 plt.tight_layout()
672 plt.savefig("img/importance_spiro_features.png")
673
674 # In[165]:
675
676 f_scores = f_score_summary(scorers)
677 print f_scores.to_latex()
678 f_scores
679
680 # ### Polynomial Combinations
681
682 # In[166]:
683
684 def create_poly_features(x_data, names):
685     # create new features base on Polynomials of the original best two predictors
686     poly = preprocessing.PolynomialFeatures(2, include_bias=False, interaction_only=True)
687     poly_features = pd.DataFrame(poly.fit_transform(x_data[names]), index=x_data.index)
688     poly_features.columns = ["POLY_%d" % i for i in poly_features.columns]
689     return poly_features
690
691 poly_features = create_poly_features(Xp, ['PRE4', 'PRE5'])
692 Xp_poly = pd.concat([Xp, poly_features], axis=1)
693 feature_importance = measure_importance(Xp_poly, Yp)
694 Xp_poly.drop(feature_importance[feature_importance == 0].index, inplace=True, axis=1)
695
696 # In[167]:
697
698 for model in models:
699     model['train_data'] = (Xp_poly, Yp)
700
701 scorers = pipeline.repeated_cross_fold_validation(models, n=10, k=5)
702
703 # In[168]:
704
705 get_ipython().magic(u'matplotlib inline')
706 for key, scorer in scorers.iteritems():
707     scorer.plot_roc_curve(mean_label=key, mean_line=True, show_all=False)
708
709 plt.plot(np.arange(0,1.1, 0.1), np.arange(0,1.1, 0.1), '—')
710 plt.savefig("img/roc_poly_features.png")
711
712 # In[169]:
713
714 feature_importance.plot(kind='barh')
715 plt.xlabel('Gini Impurity')
716 plt.tight_layout()
717 plt.savefig("img/importance_poly_features.png")
718
719 # In[170]:

```

```

728 f_scores = f_score_summary(scorers)
729 print f_scores.to_latex()
730 f_scores
731
732
733
734 # ## Resampling the Dataset
735 #
736 # Testing whether using resampling improves performance
737
738 # ### Testing with regular Over/Under sampling
739
740 # In[43]:
741
742 splitter = pipeline.OverUnderSplitter(test_size=0.2, under_sample=0.4, over_sample=0.8)
743 overunder_scorers = pipeline.monte_carlo_validation(Xp, Yp, models, splitter, n=50)
744
745
746 # In[45]:
747
748 for key, scorer in overunder_scorers.iteritems():
749     scorer.plot_roc_curve(mean_label=key, mean_line=True, show_all=False)
750
751
752 # In[46]:
753
754 f_score_summary(overunder_scorers)
755
756
757 # In[48]:
758
759 summarise_scorers(overunder_scorers)
760
761
762 # ### Testing with SMOTE + Undersampling
763
764 # In[171]:
765
766 smote_params = {'kind': 'regular', 'k':3, 'ratio': 0.8, 'verbose': 1}
767 splitter = pipeline.SMOTESplitter(test_size=0.2, under_sample=1.0, smote_params=smote_params)
768 smote_scorers = pipeline.monte_carlo_validation(Xp, Yp, models, splitter, n=50)
769
770
771 # In[172]:
772
773 for key, scorer in smote_scorers.iteritems():
774     scorer.plot_roc_curve(mean_label=key, mean_line=True, show_all=False)
775
776 plt.plot(np.arange(0,1.1, 0.1), np.arange(0,1.1, 0.1), '—')
777 plt.savefig("img/roc_smote.png")
778
779
780 # In[173]:
781
782 smote_f_scores = f_score_summary(smote_scorers)
783 print smote_f_scores.to_latex()
784 smote_f_scores
785
786
787 # ## Best Classifier
788
789 # In[199]:
790
791 spiro_features = create_spiro_features(Xp)
792 poly_features = create_poly_features(Xp, ['PRE4', 'PRE5'])
793 Xp_all = pd.concat([Xp, poly_features, spiro_features], axis=1)
794 Xp_all.drop(['DGN_1', 'DGN_8'], axis=1, inplace=True)
795 for model in models:
796     model['train_data'] = (Xp_all, Yp)
797
798
799 # In[200]:
800
801 scorers = pipeline.repeated_cross_fold_validation(models, n=10, k=5)
802
803
804 # In[201]:

```



```

805 for key, scorer in scorers.iteritems():
806     scorer.plot_roc_curve(mean_label=key, mean_line=True, show_all=False)
807
808 plt.plot(np.arange(0,1.1, 0.1), np.arange(0,1.1, 0.1), '—')
809
810
811 # In[202]:
812
813 f_scores = f_score_summary(scorers)
814 f_scores
815
816
817 # ## Prediction on Test Set
818 #
819 # Finally, based on the best combination of techniques used in the preceeding sections, and using the
820 # classifier with the best AUC performance, make probabilistic predictions based on the unlabelled test
821 # data.
822
823 # In[207]:
824
825 Xtest, _ = preprocess(test, y_data=None)
826 Xtest = Xtest.drop(['test_id'], axis=1)
827
828 test_spiro_features = create_spiro_features(Xtest)
829 test_poly_features = create_poly_features(Xtest, ['PRE4', 'PRES'])
830 Xtest = pd.concat([Xtest, test_spiro_features, test_poly_features], axis=1)
831
832 print Xtest.columns.size
833 print Xp_all.columns.size
834
835 # In[224]:
836
837 final_model = models[3]['model']
838 final_model.fit(Xp_all, Yp)
839 predicted_prob = pd.Series(final_model.predict_proba(Xtest)[: , 1])
840 predicted_label = pd.Series(final_model.predict(Xtest))
841
842
843 # In[225]:
844
845 final_submission = pd.concat([test.test_id, predicted_label, predicted_prob], axis=1)
846 final_submission.columns = ['test_id', 'predicted_label', 'predicted_output']
847 class_balance_summary(predicted_label)
848 final_submission

```

./src/analysis.py

A. ROC Analysis Module

```

1 from sklearn import metrics
2 import numpy as np
3 from scipy import interp
4 import matplotlib.pyplot as plt
5 from scipy.stats import threshold
6
7 class ROCAnalysisScorer:
8     """ Custom scorer to capute both the AUC score and the ROC curves.
9
10    A normal scorer for sklearn cannot capture multiple scores at once.
11    This object allows use to calculate multiple quantiies in one pass
12    of a cross validation object.
13
14    This will also capture the F scores
15    """
16    def __init__(self):
17        self.rates_ = []
18        self.aucs_ = []
19        self.fhalf_scores_ = []
20        self.f2scores_ = []
21        self.flscores_ = []
22
23    def __call__(self, ground_truth, predictions, **kwargs):
24        """ Custom __call__ function to make the object look like a

```

```

25     function thanks to python's duck typing.
26     """
27     return self.auc_score(ground_truth, predictions, **kwargs)
28
29 def auc_score(self, ground_truth, predictions, **kwargs):
30     """ Calculate the AUC score for this particular trial.
31
32     This will also calculate the F scores and ROC curves
33
34     Args:
35         ground_truth: vector of class labels
36         predictions: vector of predicted class labels
37
38     Returns:
39         AUC score for this trial
40     """
41
42     # calculate f scores
43     thresholded = threshold(predictions[:, 1], threshmin=0.5)
44     thresholded = threshold(thresholded, threshmax=0.5, newval=1.0).astype(int)
45     fhalf_score = metrics.fbeta_score(ground_truth.astype(int), thresholded, beta=0.5)
46     f2_score = metrics.fbeta_score(ground_truth.astype(int), thresholded, beta=2)
47     f1_score = metrics.fbeta_score(ground_truth.astype(int), thresholded, beta=1)
48
49     # calculate ROC curve and AUC
50     fpr, tpr, _ = metrics.roc_curve(ground_truth, predictions[:, 1])
51     area = metrics.auc(fpr, tpr)
52
53     self.fhalf_scores_.append(fhalf_score)
54     self.f2scores_.append(f2_score)
55     self.f1scores_.append(f1_score)
56     self.rates_.append((fpr, tpr))
57     self.aucs_.append(area)
58     return area
59
60 def mean_roc_metrics(self):
61     """ Compute the mean AUC and mean ROC curve """
62     mean_tpr = 0.0
63     mean_fpr = np.linspace(0, 1, 100)
64
65     for fpr, tpr in self.rates_:
66         mean_tpr += interp(mean_fpr, fpr, tpr)
67
68     mean_tpr = mean_tpr / len(self.rates_)
69     area = metrics.auc(mean_fpr, mean_tpr)
70     return mean_fpr, mean_tpr, area
71
72 def plot_roc_curve(self, title=None, labels=None, show_all=True, chance_line=False, mean_line=False,
73 mean_label="Mean"):
74     """ Plot ROC curves for all trials attached to this object
75
76     Args:
77         title: the title for the plot
78         labels: the labels to use for each line. Either list, string or None
79         show_all: show all plots or only the mean line
80         chance_line: show the chance line
81         mean_line: show the mean line
82         mean_label: label for the mean line
83     """
84     if show_all:
85         for i, ((fpr, tpr), area) in enumerate(zip(self.rates_, self.aucs_)):
86             if labels is None:
87                 name = "ROC %d" % (i+1)
88             elif isinstance(labels, list):
89                 name = labels[i]
90             elif isinstance(labels, str):
91                 name = labels
92             plt.plot(fpr, tpr, label='%s AUC = %0.4f' % (name, area))
93
94     if mean_line and len(self.rates_) > 1:
95         mean_fpr, mean_tpr, area = self.mean_roc_metrics()
96         plt.plot(mean_fpr, mean_tpr, label="%s, AUC: %0.4f" % (mean_label, area))
97
98     if chance_line:
99         line = np.arange(0, 1.1, 0.1)
100         plt.plot(line, line, "—", label="Chance")

```

```

101         if title is None:
102             title = 'Receiver operating characteristic'
103
104         plt.title(title)
105         plt.xlabel('False Positive Rate')
106         plt.ylabel('True Positive Rate')
107         plt.legend(loc="lower right", prop={'size': 10})
108
109     def plot_confusion_matrix(cm, title='Confusion matrix', cmap=plt.cm.Blues):
110         """ Plot a confusion matrix
111
112         Args:
113             cm: square matrix representing the confusion matrix
114             title: title to show on the plot
115             cmap: the colour map to use
116         """
117         plt.imshow(cm, interpolation='nearest', cmap=cmap)
118
119         width, height = cm.shape
120         for x in xrange(width):
121             for y in xrange(height):
122                 plt.annotate(str(cm[x][y]), xy=(y, x),
123                             horizontalalignment='center',
124                             verticalalignment='center', size=20)
125
126         plt.title(title)
127         plt.colorbar()
128         plt.show()
129

```

./src/roc_analysis.py

B. Pipeline Module

```

1 from roc_analysis import ROCAnalysisScorer
2 from sklearn import cross_validation
3 from sklearn.metrics import make_scorer
4 from unbalanced_dataset import SMOTE, UnderSampler, OverSampler
5
6
7 def cv_pipeline(model, x_data, y_data, cv=None):
8     """ Cross Validate a model pipeline
9
10    Args:
11        model: A sklearn Pipeline object to cross validate
12        x_data: Feature matrix
13        y_data: Class labels
14        cv: A predefined sklearn cross validation object
15    Returns:
16        A ROCAnalysisScorer object with the true/false positive rates and AUCs
17        for all folds
18    """
19    roc_data = ROCAnalysisScorer()
20    roc_data_scorer = make_scorer(roc_data, greater_is_better=True, needs_proba=True, average='weighted')
21    cross_validation.cross_val_score(model, x_data, y_data, cv=cv, scoring=roc_data_scorer)
22    return roc_data
23
24
25 def test_pipeline(model, x_data, y_data, x_test, y_test):
26     """ Test a model on a data
27
28    Args:
29        model: A sklearn Pipeline object to test
30        x_data: Feature matrix
31        y_data: Class labels
32    Returns:
33        A ROCAnalysisScorer object with the true/false positive rates and AUCs
34        for the test
35    """
36    model.fit(x_data, y_data)
37    y_hat = model.predict_proba(x_test)
38
39    test_result = ROCAnalysisScorer()
40    test_result.auc_score(y_test, y_hat)
41    return test_result

```

```

42
43
44 def score_pipeline(data, cv=None):
45     """ Score a pipeline using either cross validation (if a cross validation
46     object is provided) or using a training/test split
47
48     Args:
49         data: A dictionary object with the model and training or testing data
50         cv: Optional cross validation object to use
51     Returns:
52         A tuple of cross validation results and testing results
53     """
54
55     model = data['model']
56
57     cv_results = None
58     test_results = None
59
60     x_data, y_data = data['train_data']
61
62     if cv is not None:
63         cv_results = cv_pipeline(model, x_data, y_data, cv=cv)
64
65     if 'test_data' in data:
66         x_test, y_test = data['test_data']
67         test_results = test_pipeline(model, x_data, y_data, x_test, y_test)
68
69     return (cv_results, test_results)
70
71
72 def repeated_cross_fold_validation(models, n=10, k=5):
73     """ Run cross validation on a set of models n times
74
75     All models are tested using the same cross validation splits
76     at each iteration.
77
78     Args:
79         models: List of dictionaries containing the model
80                 and training or testing data.
81         n: number of iterations to repeat cross validation (default 10)
82         k: number of folds to use at each iteration (default 5)
83     Returns:
84         A list of scorer objects of type ROCAnalysisScorer, one for each model
85         passed.
86     """
87
88     scorers = {}
89
90     for i in range(n):
91         # create a new cross validation set for each iteration & test.
92         skf = cross_validation.StratifiedKFold(models[0]['train_data'][1], n_folds=k)
93
94         for model in models:
95             model_name = model['name']
96             if model_name not in scorers:
97                 scorers[model_name] = ROCAnalysisScorer()
98
99             results = score_pipeline(model, cv=skf)
100
101             # for each model collect the results into a single scorer.
102             # note: no average is made at this stage. The results of each
103             # of the k folds is collected into a single k * n list for
104             # the model.
105             scorers[model_name].f1scores_ += results[0].f1scores_
106             scorers[model_name].f2scores_ += results[0].f2scores_
107             scorers[model_name].fhalf_scores_ += results[0].fhalf_scores_
108             scorers[model_name].rates_ += results[0].rates_
109             scorers[model_name].aucs_ += results[0].aucs_
110
111     return scorers
112
113
114 def monte_carlo_validation(x_data, y_data, models, splitter, n=10):
115     """ Run Monte Carlo cross validation on a set of models n times.
116
117     This will randomly split the training and test data n times
118     and evaluate the performance of each model on each split.

```

```

119
120 Args:
121     x_data: Feature matrix
122     y_data: Class labels
123     models: List of dictionaries containing the model and
124             training or testing data.
125     splitter: A test splitter object that creates random training
126             test splits.
127     n: number of iterations to perform (default 10)
128 Returns:
129     A list of scorer objects of type ROCAnalysisScorer, one for each
130     model passed.
131 """
132 scorers = {}
133
134 for i in range(n):
135     x_train, y_train, x_valid, y_valid = splitter.split(x_data, y_data)
136
137     for model in models:
138         model_name = model['name']
139         if model_name not in scorers:
140             scorers[model_name] = ROCAnalysisScorer()
141
142         model['train_data'] = (x_train, y_train)
143         model['test_data'] = (x_valid, y_valid)
144
145         results = score_pipeline(model)
146
147         # for each model collect the results into a single scorer.
148         # note: no average is made at this stage. The results of each
149         # of the k folds is collected into a single k * n list for
150         # the model.
151         scorers[model_name].f1scores_ += results[1].f1scores_
152         scorers[model_name].f2scores_ += results[1].f2scores_
153         scorers[model_name].fhalf_scores_ += results[1].fhalf_scores_
154         scorers[model_name].rates_ += results[1].rates_
155         scorers[model_name].aucs_ += results[1].aucs_
156
157     return scorers
158
159
160 class TestSplitter(object):
161     """ TestSplitter base class.
162
163     This splits a feature matrix and class labels vector
164     into a training and testing split. This can be used to
165     create more complicated splitters for resampling.
166     """
167     def __init__(self, test_size=0.2):
168         self._test_size = test_size
169
170     def split(self, x_data, y_data):
171         Xtest, Xvalid = cross_validation.train_test_split(x_data, test_size=self._test_size, stratify=
y_data)
172         Ytest, Yvalid = y_data.loc[Xtest.index], y_data.loc[Xvalid.index]
173         return Xtest, Ytest, Xvalid, Yvalid
174
175
176 class SMOTESplitter(TestSplitter):
177     """ Test splitter for SMOTE datasets.
178
179     This splitter will apply smote to the training portion of the dataset
180     but will leave the testing part of the split untouched.
181     """
182     def __init__(self, under_sample=1.0, smote_params={}, **kwargs):
183         super(SMOTESplitter, self).__init__(**kwargs)
184         self._under_sample = under_sample
185         self._smote_params = smote_params
186
187     def split(self, x_data, y_data):
188         Xt, Yt, Xv, Yv = super(SMOTESplitter, self).split(x_data, y_data)
189         Xt_smote, Yt_smote = SMOTE(**self._smote_params).fit_transform(Xt.as_matrix(), Yt.as_matrix())
190         Xt_smote, Yt_smote = UnderSampler(ratio=self._under_sample).fit_transform(Xt_smote, Yt_smote)
191         return Xt_smote, Yt_smote, Xv, Yv
192
193
194 class OverUnderSplitter(TestSplitter):

```

```

195     """ Test splitter for under and/or over sampling datasets .
196
197     This splitter will apply under and/or over sampling the the training
198     portion of the dataset but will leave the testing part of the split
199     untouched.
200     """
201     def __init__(self , under_sample=1.0 , over_sample=1.0 , **kwargs):
202         super(OverUnderSplitter , self).__init__(**kwargs)
203         self._under_sample = under_sample
204         self._over_sample = over_sample
205
206     def split(self , x_data , y_data):
207         Xt , Yt , Xv , Yv = super(OverUnderSplitter , self).split(x_data , y_data)
208         Xt_smote , Yt_smote = OverSampler(ratio=self._over_sample).fit_transform(Xt.as_matrix() , Yt.
209         as_matrix())
210         Xt_smote , Yt_smote = UnderSampler(ratio=self._under_sample).fit_transform(Xt_smote , Yt_smote)
211         return Xt_smote , Yt_smote , Xv , Yv

```

../src/pipeline.py

APPENDIX D FINAL PREDICTIONS

This listing shows the final submission CSV file generated from the code in B.

```

1 test_id , predicted_label , predicted_output
2 1,0,0.0959909249466
3 2,0,0.366485835856
4 3,1,0.761545623495
5 4,1,0.709843850303
6 5,1,0.920849883091
7 6,0,0.242149793728
8 7,0,0.276558435347
9 8,0,0.351840295428
10 9,0,0.152283795503
11 10,1,0.985394679179
12 11,0,0.147004590471
13 12,0,0.336270685043
14 13,1,0.920825189002
15 14,1,0.657822484527
16 15,0,0.468365713819
17 16,0,0.442055739259
18 17,0,0.495571800452
19 18,0,0.218808918617
20 19,0,0.186103187944
21 20,1,0.683037960349
22 21,1,0.603136057345
23 22,1,0.843943753928
24 23,1,0.785502120221
25 24,1,0.562094197743
26 25,0,0.0349681280673
27 26,0,0.126758488774
28 27,1,0.837178099017
29 28,0,0.145145634277
30 29,0,0.285236418377
31 30,0,0.312892212194
32 31,1,0.60578789889
33 32,0,0.20121457422
34 33,0,0.20394148876
35 34,1,0.693834215385
36 35,0,0.225362734386
37 36,0,0.067643429401
38 37,0,0.336266163162
39 38,0,0.159134323265
40 39,0,0.0983555902471
41 40,1,0.698861221197
42 41,0,0.153279819097
43 42,0,0.194327978729
44 43,1,0.884207741975
45 44,0,0.267555619336
46 45,0,0.429229065178
47 46,0,0.161159509318
48 47,0,0.149172093008
49 48,0,0.298975471662
50 49,0,0.25984646789
51 50,1,0.816215790903

```



```
52 51,0,0.140296307944
53 52,1,0.975468005433
54 53,0,0.168144520453
55 54,0,0.229665238957
56 55,0,0.454043791869
57 56,1,0.912383136769
58 57,0,0.320184003255
59 58,0,0.178311957244
60 59,1,0.860341648052
61 60,1,0.518934019112
62 61,0,0.0833337588317
63 62,0,0.130123316794
64 63,1,0.982299567278
65 64,1,0.886431415143
66 65,0,0.396071759913
67 66,0,0.304953227028
68 67,0,0.0550457347313
69 68,1,0.578768136329
70 69,0,0.446981985094
71 70,1,0.61677202839
72 71,0,0.0660299381267
73 72,0,0.0526109920857
74 73,0,0.195575018274
75 74,0,0.113239703892
76 75,0,0.401289125219
77 76,0,0.149126512876
78 77,1,0.604419014644
79 78,1,0.625668159004
80 79,0,0.114890600558
81 80,0,0.266926485031
82 81,0,0.116354916444
83 82,0,0.202088422189
84 83,0,0.0657000565686
85 84,1,0.745119950263
86 85,1,0.864374587005
87 86,1,0.568390450972
88 87,0,0.152482082904
89 88,0,0.237317012369
90 89,1,0.73032801617
91 90,0,0.204900200812
92 91,0,0.462046573244
93 92,0,0.0897976729124
94 93,0,0.069122795763
95 94,0,0.375534261511
96 95,0,0.129356461134
97 96,1,0.952894216494
98 97,0,0.202907473556
99 98,0,0.0842044147829
100 99,0,0.118616675581
101 100,1,0.773608257487
```

```
../data/submission.csv
```