

Knowledge distillation from random forests

Juan José Contreras González
School of Computer Science and
Electronic Engineering
University of Essex
Colchester, United Kingdom
jc18548@essex.ac.uk

Abstract — This report documents the development of a system that will compress a Random Forest Classifier model into a single Decision Tree Classifier through a knowledge distillation process. The experiments conducted will be binary classification problems, and the evaluation of the compression will be by comparing the performance of the decision tree trained through the distillation process with a tree trained with the original dataset. The distillation process consists on training a Random Forest Classifier with the dataset and then transferring the knowledge to a single decision tree by the generation of a new dataset for multi-class classification that includes the probabilities, as calculated by the Random Forest Classifier. The new dataset fits the probabilities accordingly through a binning process, and when the model predicts, the bins are divided again to form the original classes from the Random Forest Classifier.

The outcome of this work will contribute to the assessment of model compression in Random Forest Classifiers for binary classification problems.

Keywords—Random Forests, Decision Trees, Knowledge Distillation, Model Compression, Data Science

I. INTRODUCTION

A. Handling data

Data has become the central part of most of human activities. Years ago, only big companies handled data, but now, with the existence of personal computers and electronic devices, as well as the spread of wireless communication, almost every individual handles data. Every time a product is bought, a web page is visited, or something is post on social media, data is being generated and shared.

Having this in mind, analysing data and extracting important information from it has become a major discipline, because strategies can be designed in order to make the most out of the knowledge acquired from it. Client behaviour, disease diagnosis, and natural disaster prevention are just examples of applications for this knowledge. To perform this task, patterns within the data must be found.

These patterns can be found with a computer, and an algorithm is required, which is a set of instructions that must be carried out to transform an input into an output. However, there are tasks in which an algorithm cannot be applied, because both input and output are known, but the procedure to transform one into the other is unknown. Determining spam emails from legitimate messages is a good example of this, as the definition of spam messages varies from individual to individual.

When this type of information is unknown, data becomes very important, as what is lacking in knowledge can be made up for in data. Thousands of examples can be compiled, some of which are known to be considered true and some

others false to the condition that is considered (in the case of spam email, this means spam or not spam), and the goal is to learn what constitutes a certain outcome.

The process might not be identified completely, but a good and useful approximation can be constructed, and certain patterns or regularities can be detected. This is the main core of Machine Learning.

Machine Learning is programming computers to optimise performance criterion using example data or past experience. The model can be *predictive*, which means it can foresee the future, or *descriptive*, to gain knowledge from data, or a combination of both.

Since the core task of Machine Learning is making an inference from a sample, mathematical models are built through statistics theory, and computers play a two-stage role: first, when training, an algorithm is needed to solve the optimisation problem, and to store and process the existent big amounts of data. Second, once a model is learned, an efficient representation and algorithmic solution for the inference is needed.

When Machine Learning methods are applied to large databases, it is called Data Mining, and its applications are abundant, as it can be applied in finance to use in fraud detection or predictions for the stock market. However, Machine Learning is not only based on databases problems, it is also an important part of artificial intelligence, as for a system to be intelligent, it must have the ability of learning when interacting with a changing environment.

There are different types of learning:

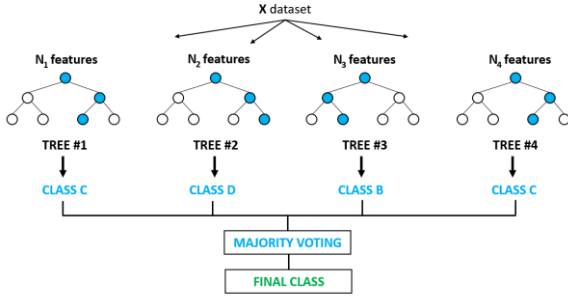
- **Supervised Learning:** it is performed when there is an input and an output, and the task is to learn the mapping from input to output, but the correct values for this output are known and provided by a supervisor.
- **Unsupervised Learning:** similar to Supervised Learning, but there is no supervisor and there is only input data, and the aim is to find regularities in this input
- **Reinforced Learning:** the output of the system is a sequence of actions, and a single one is not important; the important element is the policy that allows the goal to be reached. The program should asses the goodness of policies and learn from past action sequences to be able to generate a policy. [2]

Among Supervised Learning, there are two prediction tasks that can be performed, *Classification* and *Regression*. In classification, the objective is to learn the boundaries separating the instances of one class from the instances of all other classes. This can lead to a separation in just two classes, and is known as binary classification, or in many classes.

Regression, unlike classification, does not provide a Boolean output (the instance belongs to a class or not), but a numeric value. The learning objective in this case is not a class but a numeric function. [8]

B. Random Forests

Random Forests are a type of learning algorithm that is an ensemble of Decision Trees, under the premise that combining multiple learning models increases the overall result. Each tree depends on the values of a random vector sampled independently and with the same distribution for all trees in the forest; each one of them casts a unit vote for the most popular class at a certain input. [7]



Using Random Forests is a simple way of taking advantage of parallel computation. However, as it is an ensemble, it could take a big amount of time during the testing stage. Another throwback is the fact that if the individual models are large, as well as the dataset, the computational budget during the training stage becomes very large.

C. Distillation

In order to reduce the time required during testing, a special strategy known as “distillation” [1] can be applied. This methodology consists in using a different kind of training once a cumbersome model (in this case the Random Forest) is trained, the knowledge can be transferred to a small model that is more suitable for the deployment.

When a cumbersome model is trained to differ into many classes, the main objective is to maximise the average log probability of the correct answer, but in doing so, the model assigns probabilities to all incorrect answers. These probabilities reflect the generalisation process of the model, and while it is usually accepted that the objective function used for training reflects the true objective, models should not be trained to optimise performance on the training data, but to generalise well to new data. However, training a model to generalise well requires information, and it is not always available.

If a cumbersome model is trained and generalises well, which can happen because it produces an average of a large ensemble of different models, a small model that is trained to generalise in the same way, will perform better on test data than a small model that is trained in a normal way on the same training set that was used to train the ensemble.

An important advantage of transferring knowledge to a smaller model is that a small model can usually be trained on much less data, and with a much higher learning rate than the original cumbersome model. The approach presented for

transferring the generalisation ability from the cumbersome model to the smaller model is by using the probabilities obtained from the cumbersome model as soft targets for the other.

II. BACKGROUND

The process of transferring information from a cumbersome model to a smaller model has been performed by [1] in speech recognition, through Deep Neural Networks. In their experiments, a model used by Android voice search was set as baseline, 10 separate models of prediction and a distilled single model were compared.

The 10 separate models used the exact same architecture as the baseline and were initialised with random initial parameter values. The averaged predictions of the 10 models were used to create the soft targets.

[5] conducted a set of experiments of compressed modelling, in which they compared the performance of the compressed models with the performance of the target ensemble selection models on eight test problems. In their experiments, they compressed the ensemble selection models and MUNGE data.

An important outcome from their research is the fact that with compression, the size of the training set can be made arbitrarily large, to avoid overfitting, and reliably measure the effect of the model’s size on performance.

The distillation process is an adaptation of the training of a simple model, because it implies transferring knowledge that was acquired by a cumbersome model into a different and smaller model. This model can also be trained on less data than the bigger one and is easier to deploy.

III. METHODOLOGY

The objective of this work is to implement and evaluate the performance of a Distilled Random Forest Classifier in binary classification problems, in a similar way to the experimented conducted by [6], although they were conducted with Neural Networks.

The process makes use of a Random Forest classifier, and transfers the knowledge acquired to a single tree. This was done through the Python library scikit-learn and with three different datasets. The knowledge was transferred into the smaller model with a distillation process by using the probabilities obtained from the cumbersome model as soft targets.

Once the small models were trained by distillation, their performance was compared with the performance of similar models trained in a traditional fashion in order to evaluate the convenience of their implementation.

A. Datasets

The three chosen datasets are the following:

- **Geographical Origin of Music** – Contains audio features extracted from 1059 files, created by Fang Zhou at the University of Nottingham, Ningbo, China. [3] The objective is to predict the geographical origin of the music contained. The music is traditional, ethnic or under the ‘world’ label, as

classified by their publishers. No Western music is included, because of its important influence in music from all over the world, and would interfere with the search for aspects of music that most influence location, and the central objective is to specify a location with strong influence on the music.

The country of origin was determined by the artists' main country or area of residence. The position of the country's capital city was taken as the absolute point of origin, by latitude and longitude. The audio features from the files were extracted with the program MARSYAS, which yields a set of 68 attributes inside of a single vector. Information about chromatic pitches is also included in a different file, with a set of 116 attributes per track. All features were transformed to have a mean of 0 and a standard deviation of 1.

The dataset is contained in two files, where each file corresponds to a set of features. Both files contain the features for 1059 examples. One of them contains the audio features in the first 68 columns, and the last two correspond to the latitude and longitude of the place of origin of the music, while the other has 116 columns for features and two for the coordinates of the origin of the music.

The histogram for the distribution of this dataset is presented next:

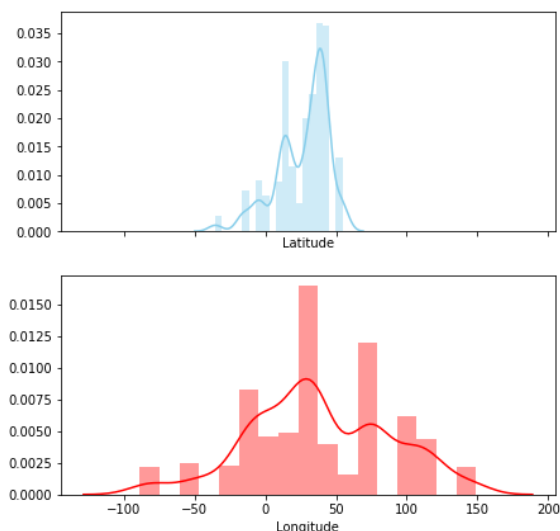


Figure 1. Data distribution for Geographical Origin of Music dataset

The statistic information for the latitude is the following:

Mean: 26.65

Median: 33.66

Var: 340.428

Std: 50.396

For the longitude:

Mean: 38.405

Median: 32.83

Var: 2539.754

Std: 50.396

- **Year Prediction MSD** – A prediction of the release year of a song from audio features. Most of the songs are Western, commercial tracks from years between 1922 – 2011, with a peak in the year 2000. It is a subset of the Million Song Dataset, a collaboration between LabROSA and The Echo Nest, prepared by T. Bertin-Mahieux. [4]

The dataset consists of 515345 examples and 90 attributes. The training and test split must allocate the first 463,715 examples for training and the last 51,630 for testing, as this division ensures that no song from a given artist ends up in both sets.

From the 90 attributes, 12 correspond to timbre average and the remaining 78 to timbre covariance. The first value in the dataset is the release year of the songs (target), ranging from 1922 to 2011.

The distribution of the data, corresponding to the year of the songs can be seen in the following histogram.

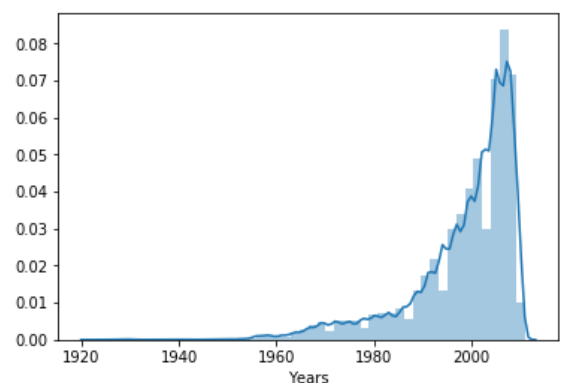


Figure 2. Data distribution for the Year Prediction MSD dataset

The statistic information for this dataset is:

Mean: 1998.397

Median: 2002

Var: 119.487

Std: 10.931

- **Default of credit card clients** – The dataset was created by I-Cheng Yeh as part of a research aimed at the case of customers' default payments in Taiwan. [5] The research employed a default payment as a binary response variable, and consists of 23 attributes, containing information regarding amount of given credit, gender, education, marital status, age, history of past payment, amount of bill statement, and amount of previous statement. The labels are given as 1 = Yes and 0 = No for the default payment.

An analysis of the relationship between the set of attributes via scatter plots, as well as histograms for the distribution of each single variable can be seen in the following chart.

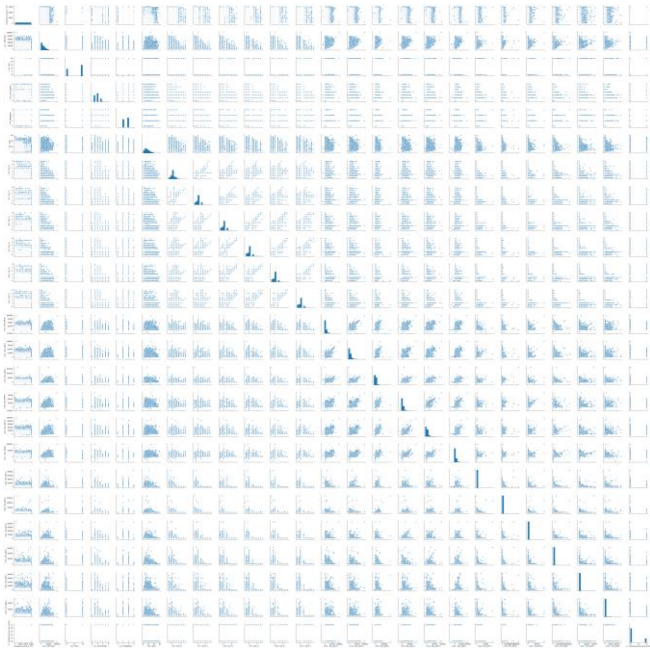


Figure 3. Data distribution and relationships for the Default of credit card clients dataset

The were revised, to avoid having any anomalies in the data, and performing any pre-processing needed, such as normalisation or deleting empty rows.

B. Methodology

One Random Forest classifier was trained for each dataset. Each of the three were constituted by a minimum of 100 trees.

Once the training was complete and the probabilities for each class were obtained, a new dataset was created, including the probabilities calculated by the forest, so that it is a multi-class classification dataset with new classes, as shown in the figure:

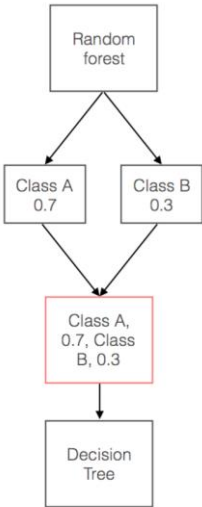


Figure 4. Knowledge transfer

Probability bins were created as binning method for the data. This allows the outputs coming from the Random Forest Classifier to be sorted in different classes, so that the Decision Tree Classifier learned to classify in more than two classes.

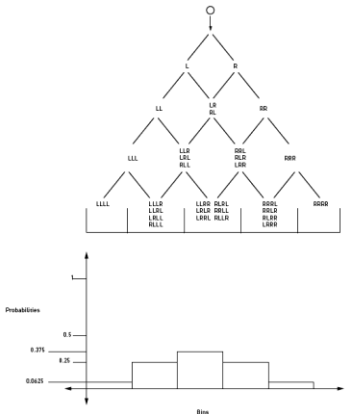


Figure 5. Probability bins

Once the tree learned and predicted examples into these bins, the information was divided again in two classes. This was performed considering that every value from the bins that is less or equal to 0.5 will belong to one class, and above 0.5 to the other class.

The probability bins are shown in the following histograms. The blue colour corresponds to class A, and the orange colour to class B. There were 10 bins, going from 0 to 9.

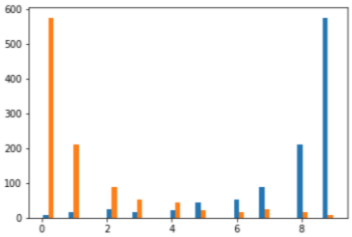


Figure 6. Geographical origin

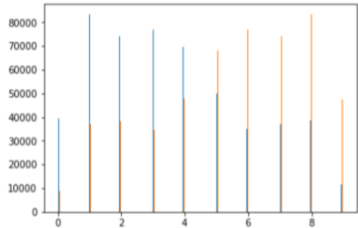


Figure 7. Year Prediction

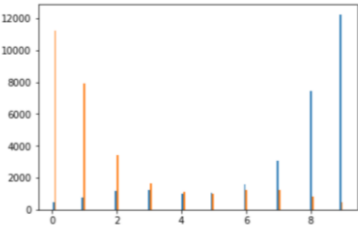


Figure 8. Credit Card

The decision tree classifiers were trained and tested with these new datasets, in order to assess their accuracy on the original datasets. The outcome was compared to the original model in terms of size, shape and performance.

C. Hyperparameters tuning

The random forest will be generated, trained and tested through the scikit learn Python library. This library allows the splitting of data into training and validation, as well as performing cross validation multiple times in order to obtain the best hyperparameters for the model.

To tune the hyperparameters of the Random Forest Classifier, a Cross Validation Grid Search will be conducted, which will perform an exhaustive search over specified parameter values for an estimator. Fitting this grid search will return a set of best parameters for the model.

This process must be performed for each of the three datasets, as the same hyperparameters will most likely not give the best result for two different cases.

Both, the random forest classifier and the decision tree classifiers will go through this hyperparameter tuning process.

IV. EXPERIMENTS

Three experiments were performed in this work, one for each of the abovementioned datasets. The main goal for each of the experiments was to be able to correctly classify examples into two classes through the distillation approach for a Random Forest Classifier, so the datasets were adjusted so that all the information can be fixed in two classes.

A. Results from previous studies

The results obtained by [1] from the experiments on speech recognition were based on Test Frame Accuracy and Word Error Rate. They are shown in the following table.

System	Test Frame Accuracy	WER
Baseline	58.9%	10.9%
10xEnsemble	61.1%	10.7%
Distilled Single Model	60.8%	10.7%

Table 1. DNN Results

It can be seen that the distilled single model and the 10-model ensemble give very similar results and they both outperform the baseline model. This information is crucial for the development of the work presented in this paper, as a Random Forest is a model that consists of the results of various single models, so the improvement over a baseline simple model will be present during training, but this already improved knowledge will then be subjected to a distillation process, so that the results can be further improved.

Besides, the results obtained by [5] reflected that the performance of the mimic models was as good or better than the ensemble models, but always better than the ones trained on the original data.

B. Setup

The dataset corresponding to the default of credit card clients is already in a form for binary classification, as the labels given are 'Yes' and 'No' for the default payment, so no pre-processing was needed.

On the other hand, both Year prediction MSD and Geographical origin of music datasets are in a format that yields regression problems, rather than classification. However, this was solved by fixing thresholds that divide the output values in two classes.

For the year prediction, the year 2000 was set as threshold so that if the example belonged to an earlier date, it was assigned as 'class A', otherwise, it was labelled as 'class B'.

In the case of the geographical origin, besides grouping the values in two classes, there is another characteristic of the dataset that must be adjusted, which is the fact that the labels are pairs of coordinates (latitude and longitude). The procedure was done considering only the longitude from the given coordinates of the examples, and comparing it with a fixed arbitrary location, which was London (-0.118092). If the absolute value of the difference is greater than 90, the example corresponds to class B, otherwise, it belongs to class A. The problem then became a binary classification in order to determine if the origin of the specified music example.

C. Objective

These experiments were established with the Decision Trees as binary classifiers from a multi-class input coming from the probabilities of belonging to each class from the Random Forest in order to prove that the knowledge transfer is occurring and that the trees are making use of the outcomes of the various models that constitute the forest to predict the classes of the original dataset.

V. DISCUSSION

The performance of the system was evaluated through different metrics. The first one was the accuracy of the system, because it reflects the percentage of unseen examples that it can classify correctly. However, this metric is not enough to fully evaluate the system, because the behaviour of the system is not understood completely unless a relationship between predictions and the correct values is present. This can be achieved through the inclusion of recall and precision.

The balance of recall and precision can be expressed as the harmonic mean of both values, and it is commonly referred to as F1 score.

The performance of the classifier was also assessed by the Kappa Statistic, as it states a proportion between its predictions and the predictions that would be expected from a perfect predictor.

The results for the conducted experiments are shown in the following table.

Distillation process				
Experiment	Accuracy	F1-score class A	F1-score class B	Kappa Statistic
Credit card	0.809	0.88	0.57	0.452
Year prediction	0.746	0.69	0.79	0.472
Geographical origin	0.848	0.91	0.53	0.438

Table 2. Distillation Results

Once these values were obtained, the difference between the results of the decision tree trained with the distilled information from the random forest and the results from the one trained normally were compared.

Normal process				
Experiment	Accuracy	F1-score class A	F1-score class B	Kappa Statistic
Credit card	0.808	0.88	0.57	0.448
Year prediction	0.752	0.69	0.79	0.486
Geographical origin	0.858	0.91	0.6	0.513

Table 3. Normal Training Results

VI. CONCLUSION

Common approaches for Machine Learning and Data Analysis have been used for some time and yielded good results. However, there are opportunities for new methods to be developed, as well as improvements to the most common ones. These improvements can be achieved by modifying certain steps during the implementation of the models, or by generating hybrid systems.

The distillation process is an approach that involves an adaptation process in the training of a simple model, as it is a transfer of the knowledge acquired by a cumbersome model into a smaller model, and its implementation derives promising results as these models perform better on test data than models that is trained in a normal way on the same training set that was used to train the cumbersome model. It is also important to mention that a small model can be trained on less data and at a higher learning rate than a bigger model.

The process has been used by [1] in a Deep Neural Network model, and good results were achieved, as the simpler neural network trained through distillation was easier to deploy than the original.

Therefore, evaluating the performance of the distillation process in Random Forests classifiers produces important information about optimising a decision process in real-time and fast deployment applications. This is specially taking into consideration the fact that, as stated by [1] in their experiments, a weighted average of 10 models outperforms a simple model, but is almost as good as a distilled model, and since Random Forests are outcomes of multiple models, this outperformance will be obtained during the initial knowledge

acquisition, and is expected to improve even more when transferred to the simple model.

From the experiments conducted in this work, it can be seen that the performance of a Decision Tree is not reduced significantly, and in more than one case, it remains the same, as in the case of the F1-scores. Both the accuracy and the Kappa score are very similar between the distilled Decision Tree and the Decision Tree trained by a normal procedure.

This method has proven to be at least as accurate as a normally trained Decision Tree. However, it takes a significant amount of time to train the Random Forest when working with very big datasets, and can be a limitation for the process, which was the case for this work, as not many attempts could be carried out due to the required time for each one.

In the case there is enough time and computational power to train the Random Forest and transfer the knowledge to a Decision Tree, this method will make use of the multiple model outcome of the Forest and will outperform a normally trained Decision Tree.

VII. REFERENCES

- [1] Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531.
- [2] Cowling, P. I., Powley, E. J., & Whitehouse, D. (2012). Information set monte carlo tree search. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(2), 120-143.
- [3] Fang Zhou, Claire Q and Ross. D. King. Predicting the Geographical Origin of Music, ICDM, 2014
- [4] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The Million Song Dataset. In *Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR 2011)*, 2011.
- [5] Yeh, I. C., & Lien, C. H. (2009). The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert Systems with Applications*, 36(2), 2473-2480.
- [6] Buciluă, Cristian, Rich Caruana, and Alexandru Niculescu-Mizil. "Model compression." *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2006.
- [7] Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
- [8] Alpaydin, Ethem. *Machine learning: the new AI*. MIT press, 2016.