# Extensive Huffman-tree-based Neural Network for the Imbalanced Dataset and Its Application in Accent Recognition

Jeremy Merrill<sup>1</sup>, Yu Liang<sup>1</sup>, Dalei Wu<sup>1</sup>
Department of Computer Science and Engineering, University of Tennessee at Chattanooga, USA
Email: nrz855@mocs.utc.edu; {yu-liang, dalei-wu}@utc.edu

Abstract—To classify the data-set featured with a large number of heavily imbalanced classes, this paper proposed an Extensive Huffman-Tree Neural Network (EHTNN), which fabricates multiple component neural network-enabled classifiers (e.g., CNN or SVM) using an extensive Huffman tree. Any given node in EHTNN can have arbitrary number of children. Compared with the Binary Huffman-Tree Neural Network (BHTNN), EHTNN may have smaller tree height, involve fewer neural networks, and demonstrate more flexibility on handling data imbalance. Using a 16-class exponentially imbalanced audio data-set as the benchmark, the proposed EHTNN was strictly assessed based on the comparisons with alternative methods such as BHTNN and single-layer CNN. The experimental results demonstrated promising results about EHTNN in terms of Gini index, Entropy value, and the accuracy derived from hierarchical multiclass confusion matrix.

Index Terms—Imbalanced classification, extensive Huffman tree neural network, CNN, audio signal, spectrogram, accent recognition

#### I. INTRODUCTION

A long-standing issue for audio classification is low accuracy due to an imbalanced data set. This paper is the result of efforts to discover a better solution to audio classification of voice accent on a skewed data set. A skew data set is a data set in which the training or validation data is skewed towards a particular class or set of classes, sometimes on an exponential level [1]. Skewed data sets can produce too many negative examples vs positive examples, causing the classifier to incorrectly predict a class for a given piece of data. Skewed data sets do not perform as well as a balanced data set. Likewise, a multi-class classification problem makes the prediction accuracy even worse. Skewed data can even alter the ability of the classifier's training by insufficiently providing enough data during training to accurately fit the model. A few concepts related to oversampling or undersampling have been utilized to help circumvent this issue [2]. Random resampling or random selective sampling also has drawbacks. This technique could also affect the classifier's ability to train and fit the model appropriately as oversampling less frequent classes doesn't provide enough data variation to build a model that is trained diverse enough to recognize slight variations in data. Moreover, random undersampling can limit the capacity of the model to fit on a diverse dataset in the same sense of undersampling. Lastly, sampling does not account for the likelihood of an input data belonging to a skewed class. Sampling indicates the classes are equally weighted to the classifier [2].

In this paper, we will explore a Huffman-tree-based approach to classification by breaking down a large multi-class problem into several smaller problems. The remainder of this paper is divided into sections. Section II will discuss related work to machine learning for audio classification. Section III will describe the canonical Huffman-Tree Neural Network (HTNN). Section IV will pronounce our contribution toward the Extensive Huffman Tree Neural Network (EHTNN). The results and comparison of each model is given in section V, and a conclusion with potential future work will be denoted in Section VI.

## II. MACHINE LEARNING ENABLED AUDIO CLASSIFICATION

The topic of audio classification using machine learning is not a new concept. Plenty of research within the artificial intelligence (AI) field of computer science has been performed. Many researchers have discovered multiple neural network techniques that can be used in order to train an AI machine to accurately classify snippets of audio based on music genre, what instrument is playing, or what words an individual is speaking. Most of these audio classification techniques use a specific type of neural network called a Convolution Neural Network (CNN). This machine learning method is mostly used to identify features of a two-dimensional image. Audio can be used as the input to a CNN whenever it is converted into a spectrogram using a mathematical signal processing technique called Fourier Transform. This converts the audio signal into a frequency vs time graph that highlights the dominant frequency at any given time within an audio snippet. This process is shown in Figure 1.

This paper will cover the practical use of using CNNs to correctly classify an audio recording of a human speaking words into a microphone by nationality through identifying frequency features of the person's accent. While this technique isn't much different than music genre classification, this paper will extend on the concept by exploring how to increase the classification accuracy of the network where there are a higher than usual number of categorical classes. Having a high number of classes reduces the ability for a single CNN to accurately classify an input. Additionally, the data set being used is unbalanced on an exponential scale. We will explore a technique for increasing accuracy of the model, reducing the time complexity, and minimizing the entropy and impurity (Gini) index.



Fig. 1. Diagram illustrating audio files as CNN input (1s sample of WAV).

Using machine learning to analyze audio is not a new concept. Several different types of neural networks exist that are used to process audio depending on the goal of the model. In practice, many audio classifications are transformed into images during the data pre-processing stage of the neural network [3]. The 2D visual representation of the audio is produced by a Fourier transform which generates a frequency vs time plot. This image is useful for classification and sequential analysis because a Fourier transform highlights dominant frequencies at any given moment throughout the duration of the audio sample [3]. A Fourier transform (often referred to as a spectrogram) is then fed into an image based neural network model.

#### A. Data Preparation for Accent Recognition

Many datasets are available for this type of work. For this paper, the "common-voice" dataset was retrieved from Kaggle. This dataset contains hundreds of thousands of mp3 files containing words spoken by many different people. On average, each mp3 file is only a few seconds long. A quick sample of a few of the files reveals the person is either speaking a sentence or multiple random words, some with a brief pause in between. The dataset is accompanied by a meta data file that contains a record for each mp3 file - indicating the corresponding filename, words spoken, and some demographics about the person speaking (gender, age, and accent) [4]. To simplify the training process, a subset of the data was used in order to speed up the training and to prevent the learning engine from exceeding system resources. The first twenty-five thousand records in the training dataset that contained a value for the accent field were used.

To further prepare the data, each audio file was sampled for 1 second to create the same spatial resolution for the spectrogram. The audio recordings provided in the data set are not of equal length and the resolution mismatch would cause false positives in the neural network's convolution and max-pooling layers.

Each spectrogram image for the model input was prepared using the ImageDataGenerator class within the Keras API. This class simplifies the import of training and validation data into a tensor data type as well as loading the correct classification labels for each piece of data. Additionally, the spectrograms were resized to a square image (which is ideal for CNNs) and converted to grayscale to eliminate color differences on the dataset.

#### B. Support Vector Machine Enabled Audio Classification

Support Vector Machines (SVM) are occasionally used to classify images in a machine learning model where a finite feature can be chosen [3][5]. The SVM can be used to classify images using more than one feature; however, additional features are computationally more expensive quadratically [5].

For this paper, we chose CNN over SVM as the model for these reasons. However, SVMs are most applicable to binary classification models where the neural network must classify an image input between one class or another [5]. SVMs can be used to classify images into multiple classes; however, they must be implemented in a chained manner [5].

### C. Convolutional Neural Network Enabled Audio Classifica-

Convolutional Neural Networks (CNN) are also used for classification by using filters to identify multiple features within the image that are similar to other images within the same class [6]. Each image from the input is trained by convoluting a subset of pixels from the image with a kernel matrix, resulting in a convoluted matrix. Weights are then applied to the parameters of the network during training as the network is identifying important features and enhancing the prediction accuracy of the network [6]. The convolutional layer is followed by a pooling layer that extracts the most important features from the convolution and reduces the feature dimensions [6]. Finally, an activation function is used to determine which neurons of the neural network are important to the classification. Recent studies have concluded that the Rectified Linear Unit (ReLU) function is best for CNNs because it does not activate all neurons at once [6]. ReLU is also been shown to have better performance than other functions, such as sigmoid [6].

#### III. BINARY HUFFMAN TREE NEURAL NETWORK

As a remedy to the skew data set, a binary Huffman tree neural network was proposed in [7] as a refined method to increase classification and prediction accuracy. However, due to the computational limitations of SVM-based neural networks for data with multiple features, we are proposing a CNN-based Huffman Tree Neural Network.

#### **Algorithm 1:** Construction of binary Huffman tree

Result: The binary Huffman tree for BFTNN

- 1 initialization:  $L = \{\langle accent_i, frequency_i \rangle\}_{i=1}^m$ ;
- 2 while (|L| > 1) do
- 3 Sort all the accents in *L* in descending order of the frequencies;
- Insert first two accents with smallest frequency into the tree;
- 5 Combine the above two accents into a joint one;

6 end

#### A. Huffman Tree Architecture and Numerical Analysis

The Huffman tree increases accuracy by predicting if inputs are classified as the highest frequently occurring class first, followed by another classification of the next higher occurring class, and so on [7][8][9]. As a result, this technique requires a numerical analysis to be performed on the data set. The analysis will determine how to initially configure the Huffman tree. Each class is listed in descending order with the frequency of occurrence within the dataset. Such an analysis is given in Table I.

TABLE I FREQUENCY OF COMMON VOICE ACCENTS

Code	Accent	Frequency
a	us	11974
b	england	5751
c	india	1806
d	australia	1677
e	canada	1486
f	scotland	565
g	africa	445
h	newzealand	422
i	ireland	364
j	philippines	125
k	wales	104
1	bermuda	80
m	malaysia	75
n	singapore	50
o	hongkong	42
p	southatlantic	34

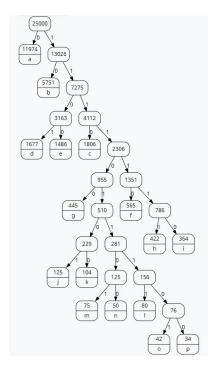


Fig. 2. Binary Huffman Tree Neural Network (BHTNN) for the Common Voice Accents

Huffman coding defines the algorithm and procedure in order to build a visual binary tree that represents the frequency table. Starting with the classes of the lowest frequency, create two tree nodes. Create a parent node with the sum of the frequencies of the child nodes. As you move up the table, if the next class's frequency exceeds the total of the recently created parent node, a new tree should be created and then combined with the current tree with a new parent node [7][8][9]. Figure 2 shows the resulting binary Huffman Tree derived from the Common Voice Accents data set. The tree nodes represent binary neural network classifier.

#### B. BHTNN Based on Support Vector Machines

Since SVMs are binary classification networks; they must be chained together when implementing a multi-class output



Fig. 3. Illustration of CNN Model

[5]. Such a method can be achieved by implementing a Binary Huffman Tree Neural Network (BHTNN) with multiple individually trained SVM networks. Using a binary tree allows the machine learning engineer to create multiple SVM networks in a one-against-all or one-against-one dataset. Eventually, each class will be represented within the decision tree [5]. Using a BHTNN tree will increase the classification accuracy because the most commonly occurring class will be the first neural network in the model in a one-against-all implementation [7]. Using the Huffman encoding algorithm, a binary tree is developed with all childless nodes representing a class and all parent nodes referring to an individually trained neural network [7].

#### C. BHTNN Based on CNN

Each decision point within the BHTNN is a separately trained CNN. Each network is configured for binary classification as a one-against-all, some-against-all, or a one-against-one data-set. At the top of the tree, the first classification is a one-against-all. Mid-way through the tree, there are a few some-against-some networks. It should be noted that these networks are no more than three classes-against-all. In this case, the network is predicting if an input is any of the two or three classes vs everything else. Accuracy suffers slightly in the some-against-all and the some-against-some networks. The one-against-one networks are simpler than the others as they are only predicting between two classes.

The model is created on a Python Anaconda interpreter using TensorFlow's Keras API library for machine learning. The CNN contains five convolution layers each followed by a max-pooling layer. These convolution and max-pooling layers highlight and pinpoint the identifiable features of each input. The model is then followed up by a dropout layer that randomly eliminates some of the input units to prevent overfitting. The dropout factor is set to 0.5 due to the vast nature of multiple features within the data on spectrograms. The model then passes through a flatten layer and two dense layers. The final dense layer's output size is one since the loss objective function is binary-crossentropy (indicating a yes or no which is illustrated by the 1 or 0 on the tree path in Figure 2). Up to forty training epochs were performed for each network. For the first two CNNs, a steps-per-epoch value was specified to prevent over-fitting due to the high number of training data for the two most frequent classes. The dropout factor was reduced to 0.2 to compensate for the reduced training steps. The CNN at each decision point in the BHTNN is illustrated in Figure 3

A few other model callbacks were used to further prevent over-fitting. An early stopping callback was established to stop the training if the validation accuracy had not increased after fifteen epochs. Also, the ReduceLRonPlateau function was utilized – which reduces the learning rate of the optimizer function if the validation accuracy does not increase after five epochs. Finally, a model checkpoint was used to save the parameter weights of the epoch that yields the highest validation accuracy. These peak accuracy values are given in the results section.

#### D. Binary Huffman Tree Neural Network

Huffman encoding is typically used as a lossless compression algorithm by indicating frequently used values first with less frequently used value later. The bit-storage technique is optimized by using less bits for commonly occurring values and more bits for less commonly occurring values. Visually, this would represent a higher position in a binary tree, and a lower position in a binary tree respectively. The bit length indicates how many branches to travel in the tree in order to find a particular value. More bits results in a higher time complexity for the worst case scenario [7].

This paper explores a machine learning model that will increase the accuracy of a traditional CNN in order to achieve a better performing model. These models are typically ideal for image classification; however, when CNNs have a higher number of output classes, the prediction accuracy of the tree decreases [6]. This technique will extend the CNN by encompassing multiple CNNs in a Huffman Tree Architecture (HTA). The goal of this method is to improve accuracy of a multi-class CNN by implementing a HTA with binary CNNs. This will improve the classification accuracy of an exponentially unbalanced dataset (as compared to a flat multiclass categorical cross-entropy CNN). HTAs within neural networks have been used before, but with SVMs [8][9]. Since SVMs are not computationally efficient for multiple features within an input image, this paper will explore a CNN-based HTA, an extensive CNN-based HTA with a reduced height, and an improved CNN-based HTA, utilizing entropy and impurity (Gini) index calculations to re-arrange nodes for better accuracy.

#### IV. EXTENSIVE HUFFMAN TREE NEURAL NETWORK

For a given n-class classification problem, Binary Huffman tree neural network (BHTNN) always consists of (at most) (n-1) component binary classification neural networks, therefore BHTNN lacks of flexibility in architecture configuration. In the worst case (the height of the tree is n-1), a decision may be made by passing through n-1 neural networks. As a remedy, we propose an Extensive Huffman-Tree Neural Network (EHTNN):

- Each non-leaf EHTNN node has an arbitrary number of (two or more than two) sub-trees. As a result, a lesser number of component neural networks (binary or multiclass) are needed. The height of the resulting EHTNN can be reduced as well.
- The EHTNN is improved according to specific measurement metrics such as entropy or Gini index.

Different from BHTNN (constructed using Algorithm 1), EHTNN fabricates multiple neural networks (or nodes) together using an extensive Huffman tree, whose construction is described by Algorithm 2.

```
Algorithm 2: Construction of extensive Huffman tree
```

```
Result: The extensive Huffman tree for EHTNN
1 initialization: L = {⟨accent<sub>i</sub>, frequency<sub>i</sub>⟩}<sub>i=1</sub><sup>m</sup>;
2 while (| L |> 1) do
3 | Sort accents L = {⟨accent<sub>i</sub>, frequency<sub>i</sub>⟩}<sub>i=1</sub><sup>m</sup> in ascending order of the frequencies;
4 | Compute the maximal k such that std({frequency<sub>i</sub>⟩}<sub>i=1</sub><sup>k</sup>) ≤ α * frequency<sub>0</sub>;
5 | /* frequency<sub>0</sub> is the smallest frequency */;
6 | /* std stands for standard deviation */;
7 | Insert first k (k ≥ 2) accents with smallest frequency into the tree;
8 | Combine the above k accents into a joint one;
9 end
```

In order to select the appropriate numbers of classes to move to the canonical CNN portion of the EHTNN, algorithm 2 is considered. The goal of the algorithm is to find the ideal number k of classes that will result in the standard deviation of k classes being lower than some arbitrary fraction,  $\alpha$  of the frequency of the smallest class.

Accuracy of each decision point in the tree is one metric to consider when determining the overall effectiveness of a BHTNN with an unbalanced data-set. However, there is a concern for point-of-no-return if an inaccurate decision is made at a higher point in the tree for an input that belongs to another class. Further learning on an input data that does not belong to either class in a BHTNN can negatively affect the overall performance of each neural network in the tree. The EHTNN is a concept to combine a multi-class CNN with a Huffman tree to converge the benefits of both models. Referring back to the frequencies given from the numerical analysis in Table I, the unbalance of the data becomes less apparent as we move down the table. This is the point in which we stop the Huffman tree and convert to a multi-class CNN to handle the remaining classes that are more balanced to each other as compared to the more frequently occurring classes at the top of the table. For our analysis, the splitting point will be after the fifth class in the table, resulting in a tree as illustrated in Figure 4. Additionally, reducing the height of the tree also decreases the time complexity and training efforts for the entire EHTNN.

#### A. Measurement Metrics of the EHTNN

Decision-tree-based neural networks can be measured for efficiency using multiple parameters. Two major factors are the entropy value and the impurity (Gini) index. By applying these concepts to the EHTNN, we can improve the configuration of tree network through minimizing the Gini index and the entropy.

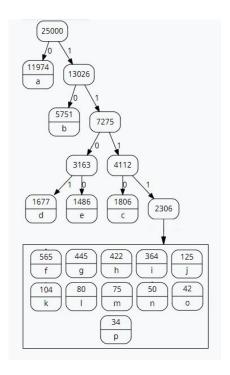


Fig. 4. A sample Extensive Huffman Tree Neural Network (EHTNN) of the Common Voice Accents

a) Entropy: The information entropy value is a statistical measurement of the uncertainty in an outcome. Values are between zero and one. The goal of the decision tree is to arrange the decisions in the tree to minimize the entropy value. Our EHTNN is rearranged by swapping the location of the US and England classes. The entropy value is then re-computed to determine the improvement of the decision tree arrangement. The entropy formula is given by equation 3 where  $p_j$  represents the probability of the class within its respective network. Each class within the tree has a computed entropy and it is weighted by multiplying by the frequency of each class within the dataset, given by equation 1. The weighted entropy values are summed to determine the final weighted entropy of the tree [10]:

$$W\mathcal{H} = \sum_{m} w_m * \mathcal{H}_m, \tag{1}$$

where the weight  $w_m$  of m-th component classifier is defined as

$$w_m = \frac{frequency_m}{\sum_i frequency_i} \tag{2}$$

and the Entropy of m-th component classifier is

$$\mathcal{H}_m = -\sum_j p_j \log_2 p_j. \tag{3}$$

where  $p_i$  is the probability of class j.

b) Gini Index: The Gini index is a measurement that indicates a level of dataset contamination or impurity. Similar to information entropy, the Gini index is a value between zero and one and a higher value indicates more data contamination. As a result, the goal of a decision tree is to reduce the value of the Gini index. A lower measurement value would signify the purity of the dataset. As a continuation of the

TABLE II

CONFUSION MATRIX ABOUT THE FIRST COMPONENT CLASSIFIER FOR US VS NOTUS IN EHTNN (BASED ON 67% ACCURACY)

		Actual Class	
		US	NotUS
Predicted	US	8022	3952
Class	NotUS	4299	8727

TABLE III

CONFUSION MATRIX ABOUT THE SECOND COMPONENT CLASSIFIER FOR ENGLAND VS NOTENGLAND IN EHTNN (BASED ON 72% ACCURACY)

		Actual Class	
		England	NotEngland
Predicted	England	4130	1621
Class	NotEngland	2051	5224

information entropy arrangement, this paper will apply a weighted Gini index (given by equation 4 to measure the impurity of extensive Huffman tree:

$$WG = \sum_{m} w_m * G_m \tag{4}$$

where the weight  $w_m$  of m-th component classifier is defined as Equation 2 and the corresponding Gini-index is defined as

$$\mathcal{G}_m = 1 - \sum_j p_j^2. \tag{5}$$

In equation 5,  $p_i$  is the probability of class j.

c) Hierarchical Multiclass Confusion Matrix: As addressed above, EHTNN consists of multiple neural network classifiers organized in a hierarchical architecture. The accuracy of *m*-th component classifier can be measured using micro-averaged F1-score. Precision and Recall, which are defined as:

$$Accuracy_m = Precision_m = Recall_m = \frac{\sum_i M_{ii}}{\sum_{i \neq i} M_{ii}}, \quad (6)$$

where M indicates the multiclass confusion matrix for component classifier m. The confusion matrices for the first two component classifiers of EHTNN are given in tables II and III. Equation 6 illustrates that when we calculating the metrics globally, all the measures (accuracy, precision, recall, and micro-average F1-score) become equal. Following equation 7, we have a global accuracy value of EHTNN that can be formulated as the weighted sum of the component classifiers.:

$$W\mathcal{A} = \sum_{m} w_{m} * Accuracy_{m}. \tag{7}$$

Matthews correlation coefficient (MCC) or phi coefficient was not discussed in this work because it is only applicable for a completely binary classification Huffman tree. Therefore, it is not applicable for EHTNN.

#### V. EXPERIMENTAL RESULTS

In order to establish a baseline, the same CNN model was trained in categorical mode; however, the loss function was changed to categorical-crossentropy with the number of output layers on the final dense layer set to sixteen (the total number

TABLE IV
RESULTS FROM EACH CNN IN THE HTA

Codes	CNN	Accuracy
a	US/notUS	0.6718
b	England/notEngland	0.718
	Australia or Canada/not	0.65625
de	Australia/Canada	0.52821
c	India/notIndia	0.625
	Scotland, NZ, or Ireland/not	0.675
f	Scotland/notScotland	0.57971
hi	Ireland/NZ	0.6842
g	Africa/notAfrica	0.60784
	Philippines or Wales/not	0.59259
jk	Philippines/Wales	0.71429
	Malaysia or Singapore/not	0.6923
mn	Malaysia/Singapore	0.625
1	Bermuda/notBermuda	0.667
op	HongKong/SouthAtlantic	0.667

of unique classes in [4]). This is the control model to make an improvement comparison to the BHTNN and EHTNN models. After several epochs, the learning converged on 45% accuracy. In comparison, the BHTNN and EHTNN both yielded (on average) a 65% accuracy. By the accuracy values alone, it can be observed that the Huffman-tree based models provide a better classification prediction for spectrograms. The validation accuracy of each CNN is given in Table IV.

As indicated, each CNN and class are not equally weighted to others, so the overall accuracy of the BHTNN and EHTNN cannot be computed by a mean. We can compute the total accuracy of each model by using a weighted mean. The formula for a weighted arithmetic mean is given in Equation 7.

Applying formulas 1, 4, and 7, we obtain the results for each method as indicated in table V.

TABLE V
RESULTS FROM EACH HTNN METHOD

Method	Accuracy	Entropy	Gini
Flat CNN	47%	0.65	0.72
Binary HTNN	66%	0.99	0.49
EHTNN ( $\alpha = 1.0$ )	66%	0.96	0.53
EHTNN ( $\alpha = 0.5$ )	70%	0.89	0.48

As indicated in the results above, the accuracy remains steady at 66% for both the BHTNN and the EHTNN, but rose slightly to 70% for the improved EHTNN. The entropy decreased from 0.96-0.99 to 0.89 on the improved EHTNN while the Gini index fell from 0.72 on the flat CNN to 0.48 for the improved EHTNN. When the EHTNN is improved for lower entropy and Gini index values, the accuracy of the model increases [10][11].

#### VI. CONCLUSION AND FUTURE WORK

The work of this paper yielded a higher accuracy value for the prediction model by using an HTNN. As stated, the Huffman tree is not a new concept for machine learning models. It was previously used with an SVM-based model. Due to the limitation of SVM-based decision trees, this paper explored increasing model prediction accuracy when a high number of output classes are present and the input must be

trained on multiple features. The CNN-based HTNNs yielded an accuracy increase of roughly twenty percent to twenty-five percent (20%-25%). While the entropy value and Gini indices decreased with a HTNN vs canonocial, the rearranged EHTNN lowered these values for improved accuracy. Even though a lower entropy and Gini index are ideal in a machine learning model, the higher values are an acceptable trade-off for higher accuracy.

Future work can be explored to further develop and optimize the EHTNN by utilizing more data preparation. Expanding on the audio processing functionality before generating a spectrogram can decrease the overfitting potential of the model. The spectrograms produced by the Python matplotlib library were adequate for use in the model; however more analysis can be performed to detect a higher audio amplitude so the data sample isn't a brief moment of silence.

Additionally, more work is needed to further arrange a decision tree based on the lowest values for entropy and the Gini index. An ideal implementation would be to compute the entropy and Gini index for each possible arrangement of the decision tree and train the networks as such. Further research and studies can be conducted to re-arrange the decision tree after initial training to keep the accuracy high while keeping the entropy and Gini index low as the tree predicts and obtains more data throughout the lifecycle of the tree.

#### ACKNOWLEDGEMENT

This work was jointly sponsored by National Science Foundation (NSF) grant numbers 1924278, 1761839 and 1647175.

#### REFERENCES

- [1] J. Cervantes, F. García-Lamont, A. López, L. Rodriguez, J. S. Ruiz Castilla, and A. Trueba, "Pso-based method for svm classification on skewed data-sets," in *Advanced Intelligent Computing Theories and Applications*, D.-S. Huang and K. Han, Eds. Cham: Springer International Publishing, 2015, pp. 79–86.
- [2] M. Koziarski, "Radial-based undersampling for imbalanced data classification," *Pattern Recognition*, vol. 102, p. 107262, 2020.
- [3] M. Esmaeilpour, P. Cardinal, and A. Lameiras Koerich, "A robust approach for securing audio classification against adversarial attacks," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 2147–2159, 2020.
- [4] Mozilla, "Common voice." Kaggle, 2017. [Online]. Available: https://doi.org/10.1145/3342555
- [5] R. Archibald and G. Fann, "Feature selection and classification of hyperspectral images with support vector machines," *IEEE Geoscience* and Remote Sensing Letters, vol. 4, no. 4, pp. 674–677, 2007.
- [6] Q. Li, W. Cai, X. Wang, Y. Zhou, D. D. Feng, and M. Chen, "Medical image classification with convolutional neural network," in 2014 13th International Conference on Control Automation Robotics Vision (ICARCV), 2014, pp. 844–848.
- [7] A. Moffat, "Huffman coding," ACM Comput. Surv., vol. 52, no. 4, Aug. 2019. [Online]. Available: https://doi.org/10.1145/3342555
- [8] F. Wu, W. Hu, and Y. Sun, "A novel multi-class fault diagnosis approach based on support vector machine of particle swarm optimization and huffman tree," in 2015 IEEE International Conference on Industrial Engineering and Engineering Management (IEEM), 2015, pp. 825–829.
- [9] G. Zhang, "Support vector machines with huffman tree architecture for multiclass classification," in *Progress in Pattern Recognition, Image Analysis and Applications*, A. Sanfeliu and M. L. Cortés, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2005, pp. 24–33.
- [10] D. Pomorski and C. Desrousseaux, "Improving performance of distributed detection networks: An entropy-based optimization," Signal Processing, vol. 81, no. 12, pp. 2479 2491, 2001.
- [11] D. Bertsimas and J. Dunn, "Optimal classification trees," *Mach Learn*, vol. 106, pp. 1039 1082, 2017. [Online]. Available: https://doi.org/10.1007/s10994-017-5633-9