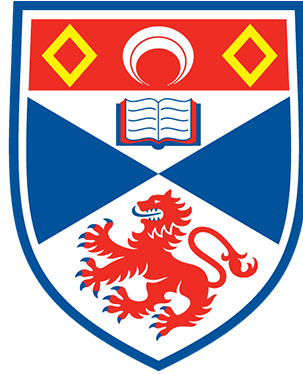


Convolutional Neural Networks  
for the Classification of  
Pseudorca crassidens Passive Acoustic Data



University of  
St Andrews

Francis Joshua Arrabaca

Supervisors: Dr Chrissy Fell and Dr Hannah Worthington

Master of Science in Data-Intensive Analysis

School of Mathematics and Statistics

University of St Andrews

August 2024

# Table of Contents

<b>List of Figures</b>	<b>3</b>
<b>List of Tables</b>	<b>6</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Pseudorca crassidens . . . . .	1
1.2 HICEAS 2017 and False Killer Whales around the Hawaiian Islands . .	2
1.3 Audio Data Collection and Curation . . . . .	4
1.4 Audio Visualisation Using Spectrograms . . . . .	5
1.5 Motivation for Study . . . . .	9
<b>2 Literature Review</b>	<b>10</b>
2.1 Machine Learning for Acoustic Data . . . . .	10
2.2 Convolutional Neural Networks . . . . .	12
2.3 Transfer Learning . . . . .	16
2.4 CNN Software Frameworks . . . . .	19
<b>3 Methodology</b>	<b>21</b>
3.1 Converting Audio WAVs to Spectrogram PNGs . . . . .	21
3.2 Transfer Learning Using Pretrained Models . . . . .	22
3.3 Addressing Imbalanced Data . . . . .	25
3.4 Hyperparameter Tuning and Final Testing . . . . .	26
3.5 Simulating Predictions with Real Data . . . . .	27
<b>4 Results</b>	<b>28</b>
<b>5 Discussion</b>	<b>35</b>
<b>6 Conclusion</b>	<b>38</b>
<b>Bibliography</b>	<b>39</b>

## Appendix A Python Code

48

# List of Figures

1.1	An illustration of an adult <i>Pseudorca crassidens</i> (False Killer Whale) by <a href="#">Dewynter (2019)</a> for the CARI'MAM project. . . . .	1
1.2	A map by <a href="#">Yano et al. (2018)</a> showing the survey area and False Killer Whale sightings. The survey includes the waters immediately surrounding the Hawaiian Islands in the lower right, but also extends further out to the PMNM area to the upper left. The black diagonal lines show the visual effort which coincide with the transects of the ship journeys. The red asterisks show visual detection only of False Killer Whales, while the blue diamonds indicate both acoustic and visual detections. . . . .	3
1.3	A simple diagram showing the relationship between frequency, wavelength and amplitude. . . . .	6
1.4	An audio file from the <i>Pelagic</i> class represented in a time domain ( <i>left</i> ) and the same file as a frequency domain ( <i>right</i> ). The left graph shows the audio amplitude of the signal across time, while the right graph shows the amplitude for the different audio frequencies. The different colours represent the audio channels in the file. . . . .	7
1.5	The audio file in Figure 1.4 now shown in a spectrogram. Unlike the previous plots, this plot shows the audio frequencies, amplitude and time in one graphical representation. While False Killer Whale whistles should be seen at around 5,000Hz and clicks at around 40,000Hz, this spectrogram shows high amplitude at 10,000Hz throughout the period, as well as around 200Hz starting around 1.0 sec. This might be indicative of other sounds such as engine and wave noises picked up by the hydrophone recorders. . . . .	8
2.1	A diagram of a simple Convolutional Neural Network by <a href="#">Phung &amp; Rhee (2019)</a> . The <i>Convolution</i> layer extracts features, while the <i>Pooling</i> layer downsamples the data. . . . .	12

- 2.2 An application of the convolution linear expression showing a convolution kernel of size  $2 \times 2$  applied on a matrix of size  $4 \times 4$  with a step size of 2-cells both horizontally and vertically. The kernel  $K$  is applied to each colour quadrant of the image  $I$ , and outputs one cell value in  $I_K$ . So  $I_K(1, 1) = (1 * 1) + (2 * 2) + (5 * 3) + (6 * 4)$  which is 44, and the process is applied to the rest of the input matrix. . . . . 13
- 2.3 A matrix of size  $4 \times 4$  is padded with extra cells around the boundaries. The recommended size of the horizontal and vertical padding  $p_h$  and  $p_w$  with respect to the kernel size  $k$  are  $p_h = k_h - 1$  and  $p_w = k_w - 1$ . Here,  $p_h$  is  $3 - 1 = 2$ , which is the same for  $p_w$ . The padded cells (*in grey*) are distributed around the perimeter and are filled with 0 for all cells (also known as *zero padding*). Here, the step size is 1 cell at a time, so the output matrix will have the same size as the input matrix. Similar to Figure 2.2,  $I_K(1, 1)$  can be computed as  $(0 * 1) + (0 * 2) + (0 * 3) + (0 * 3) + (1 * 2) + (2 * 1) + (1 * 0) + (2 * 2) + (3 * 3) = 17$ , and is similar for the rest of the cells for the output image  $I_K$ . . . . . 14
- 2.4 Illustrative examples of different pooling layers by [Yani et al. \(2019\)](#), showing how the pooling operation downsamples through the feature map. As their names suggest, Max Pooling (*left*) outputs the maximum value from each patch, while Average Pooling (*right*) outputs the average value for each patch. . . . . 14
- 2.5 An illustrative example of how a convolutional layer and pooling layer process an image of a puppy from [Pixabay \(2016\)](#). The convolutional filter identifies important features such as the puppy's eyes, nose and paws, and then downsamples the input while still keeping the important features. . . . . 15
- 2.6 A simplified diagram of transfer learning showing how a pretrained model can be used on a different dataset for a new prediction. . . . . 16

2.7	A screenshot from <i>Papers With Code: Trends</i> (2024) showing the framework usage from June 2014 to June 2024. Here, TensorFlow led the trend up until 2019 when PyTorch took a significant lead. . . . .	20
3.1	Outline of the processes done for this study. . . . .	21
3.2	Imbalanced data per batch using PyTorch’s Dataloader shows the <i>Pelagic</i> has overwhelmingly more selected instances than <i>Insular</i> . . . . .	25
4.1	Spectrograms from the first wav files in each class; <i>Insular</i> on the left and <i>Pelagic</i> on the right. The vertical streaks in the <i>Insular</i> example are not found in the <i>Pelagic</i> class. In contrast, the horizontal line in the middle of the <i>Pelagic</i> example can’t be seen in the other class. However, these may be attributed to other ambient ocean sounds. . . . .	28
4.2	The precision-recall curve ( <i>left</i> ) show how the model learns towards a precision of 0.87, which is about the same as the majority class ratio of 87.6%. The threshold plot ( <i>right</i> ) also show the effect of the severely imbalanced data, with the threshold at close to 0.0 where recall is close to 1.0. These plots show the effect of the imbalanced majority class in the dataset. . . . .	29
4.3	Using Pytorch’s <i>Weighted Data Sampler</i> , the Dataloader was able to undersample <i>Pelagic</i> instances and oversample <i>Insular</i> instances which produced a more balanced dataset across all batches. . . . .	30
4.4	ResNet Precision-Recall plots after applying <i>Weighted Data Sampler</i> . The precision and recall threshold ( <i>right</i> ) came closer to a more balanced threshold of around 0.4, however the precision-recall curve ( <i>left</i> ) while showing a balance chance level of 0.5 still showed the model line “hugging” the chance line, indicating there were still deficiencies with the model. . . . .	31
4.5	Model Losses over Epochs for best learning rate ( <i>left</i> ) and momentum ( <i>right</i> ) over 24 epochs. Despite having the best scores for hyperparameter tuning, the graphs show erratic behaviour. . . . .	32

# List of Tables

2.1	A selected list of recent studies in CNN Transfer Learning focusing on bioacoustics. While the target datasets focused on bioacoustics (bird calls, dolphin whistles, etc.), the studies all achieved superior results, even with base models that were originally trained on image data or generic audio data. . . . .	18
4.1	Validation results from transfer training using different pre-trained models, sorted by best accuracy score descending. . . . .	29
4.2	Validation accuracy results comparing the ResNet model <i>before</i> and <i>after</i> balancing the dataset. The accuracy dipped slightly after applying Pytorch's <i>Weighted Random Sampler</i> . . . . .	30
4.3	Accuracy results for manual hyperparameter tuning for Learning Rate ( <i>left</i> ) and Momentum ( <i>right</i> ). For Learning Rate, the momentum was fixed at 0.90, and subsequent Momentum tuning with the best LR of 0.01. The highlighted rows show the best scores. . . . .	32
4.4	The confusion matrix ( <i>left</i> ) showed the model prediction performance, where the True Positives represent the <i>Insular</i> class, while True Negatives are the <i>Pelagic</i> class. The model performed a decent job as it correctly identified almost two-thirds of the minority class <i>Insular</i> . . .	33
4.5	Final accuracy results comparing the best base model, the base model with weighted dataset, the weighted dataset model with the best hyperparameters, and the same model with the unseen data. Note that all model scores were tested using the validation set, <u>except</u> the final model which used the test set. . . . .	33
4.6	Predictions using a sample of the <i>Unknown</i> audio files. Most of the predictions were given with high probabilities of $>0.90$ . . . . .	34

- 5.1 Confusion matrices for (a)**EfficientNet** on the left and (b)**ResNet** on the right. Aside from the EfficientNet's poorer accuracy, it also had a Recall of 1.0 and its confusion matrix shows why this is the case; it predicted all cases as *Pelagic* (True Negatives), while none for *Insular* (True Positives). . . . . 35
- 5.2 The confusion matrices of the ResNet model (a)**before** and (b)**after** applying the *Weighted Random Sampler (WRS)*, showing how a model with lower accuracy might be more acceptable. In the left matrix, the model has a higher accuracy of 0.9256 with more True Negatives (*Pelagic*), but with fewer True Positives (*Insular*), leading to a lower precision of 0.9273. In the right matrix, the model has a slightly lower accuracy of 0.9078 due to fewer True Negatives(*Pelagic*), but more True Positives (*Insular*), leading to a much higher precision of 0.9700. . . . 36



# Code Listings

3.1	Code snippets adapted from the official PyTorch documentation by <a href="#">Chilamkurthy (2024)</a> showing how to use a CNN as fixed feature extractor in PyTorch. This example shows ResNet, but the code is similar for the other CNNs in this study. . . . .	23
3.2	Sample output when predicting for each test audio file. . . . .	27
A.1	code/a_convert_wav_to_png.py . . . . .	50
A.2	code/b_CNN_model_DenseNet.py . . . . .	52
A.3	code/b_CNN_model_efficientnet.py . . . . .	57
A.4	code/b_CNN_model_ResNet.py . . . . .	62
A.5	code/b_CNN_model_ResNet_wrs.py . . . . .	68
A.6	code/b_CNN_model_ResNet_wrs_lr.py . . . . .	75
A.7	code/b_CNN_model_ResNet_wrs_momentum.py . . . . .	80
A.8	code/b_CNN_model_ResNet_best.py . . . . .	86
A.9	code/c_make_prediction.py . . . . .	92

## **Declaration of Authorship**

I, Francis Joshua Arrabaca, hereby declare that this dissertation has been done for the degree in Msc Data-Intensive Analysis at the University of St Andrews for the academic term 2023–2024, and has not been submitted for any other degree or professional qualification.

I confirm that this dissertation has been composed by myself and was done entirely on my own, save for programming code references which have been duly acknowledged, and published work duly attributed.

## **Acknowledgements**

My deepest gratitude to my supervisors Dr Chrissy Fell and Dr Hannah Worthington of the School of Mathematics and Statistics for their extensive guidance, support and feedback. Many thanks also to the Dolphin Acoustics Vertically Integrated Project, their project head Dr Julie Oswald, and especially to Emily McCloskey for her insights in Hawaiian False Killer Whales, audio data collection and curation process.

To my wife, Mariae, and my family back home in the Philippines who have patiently supported and encouraged me from afar, I express my heartfelt love and appreciation.

Finally, to the Chevening Scholarship Program and the University of St Andrews for the opportunity to study at one of the top universities in the world, as well as to my friends, mentors and colleagues who have supported me throughout this journey, my sincerest gratitude.

Date: 13 August 2024

A handwritten signature in black ink, appearing to be 'fj' or 'fr', written in a cursive style.

## Abstract

This study involves programming experiments using Convolutional Neural Networks (CNNs) on passive acoustic data from Hawaiian False Killer Whales, and would be of interest to entry-level statisticians and computer scientists looking to study CNNs, and apply them in the field of biacoustics.

Bioacoustic classification involves manual and statistical methods of identifying animal groups using audio data. This is useful when visual or geospatial classification methods may be lacking or unfeasible. In some cases, only audio data is provided to researchers, but differences in sound may not be easily perceptible to humans, so machine learning methods are required.

This study uses bioacoustic classification to answer the question: is it possible to identify between two genetically distinct groups of False Killer Whales (*Insular* and *Pelagic*) using audio recordings alone? In answering this question, the audio data was first converted into spectrogram image data to be used in CNNs. Transfer Learning was then applied which allowed the use of pre-trained CNNs, eliminating the need to create neural networks from scratch.

The result shows that using rudimentary CNN transfer learning, it is possible to distinguish between *Insular* and *Pelagic* False Killer Whales with an accuracy of 93.89%.

**Keywords**— Bioacoustics - Convolutional Neural Network - Transfer Learning

# 1 | Introduction

This study discusses the use of Convolutional Neural Networks (CNNs) on marine mammal acoustic data, and has been written for audiences in post-graduate statistics and computer science programs. Readers are assumed to be familiar with data analysis and introductory machine learning, but no assumption is made for specific knowledge in marine bioacoustics or CNNs.

As such, statistical and machine learning concepts such as classifier models and confusion matrices are mentioned but not explained, while other concepts such as False Killer Whales, marine audio surveys, spectrograms, CNNs and Transfer Learning are discussed more thoroughly.

## 1.1 *Pseudorca crassidens*



Figure 1.1: An illustration of an adult *Pseudorca crassidens* (False Killer Whale) by [Dewynter \(2019\)](#) for the CARIMAM project.

*Pseudorca crassidens* or False Killer Whales are large oceanic dolphins that inhabit warm waters around the world, mostly in the open ocean (*Pelagic*), with some groups observed close to shore (*Insular*) ([Baird \(2009\)](#)). These animals grow to 5–6 metres in length, and are black or dark grey in colour as shown in Figure 1.1.

False Killer Whales tend to group in pods of 20 to 100 individuals ([Ridgway & Harrison \(1999\)](#)), and may associate with other marine mammals such as bottlenose dolphins. Their sounds include whistles around 5KHz lasting 0.5 seconds and echolocation clicks around 40KHz of 0.03 seconds ([DOSITS \(2021\)](#)).

In the wild, False Killer Whales are of interest as there have been unfortunate interactions with humans in the commercial and recreational fishing industries. These animals feed on the same fish species targeted by commercial fisheries (such as tunas) and have been found stealing fish from these fishing operations ([Baird \(2009\)](#)). Considered a nuisance, False Killer Whales have been killed directly or as bycatch by these fishing operators (*Pseudorca crassidens*: [Baird, R.W.: The IUCN Red List of Threatened Species 2018: e.T18596A145357488 \(2018\)](#)).

## 1.2 HICEAS 2017 and False Killer Whales around the Hawaiian Islands

The Hawaiian Islands Cetacean and Ecosystem Assessment Survey (HICEAS) 2017 was a large-scale survey of marine mammals and seabirds around the waters of the Hawaiian islands ([Yano et al. \(2018\)](#)), including the extensive waters of the Papahānaumokuākea Marine National Monument (PMNM) to the northwest of Hawaii. The survey range can be seen in Figure 1.2. The 2017 survey took place from July to December, with similar surveys conducted in 2002, 2010, and most recently in 2023 ([NOAA Fisheries \(2024\)](#)). The survey aimed to collect data in order to estimate the abundance and distribution of the animals around the Hawaiian Islands.

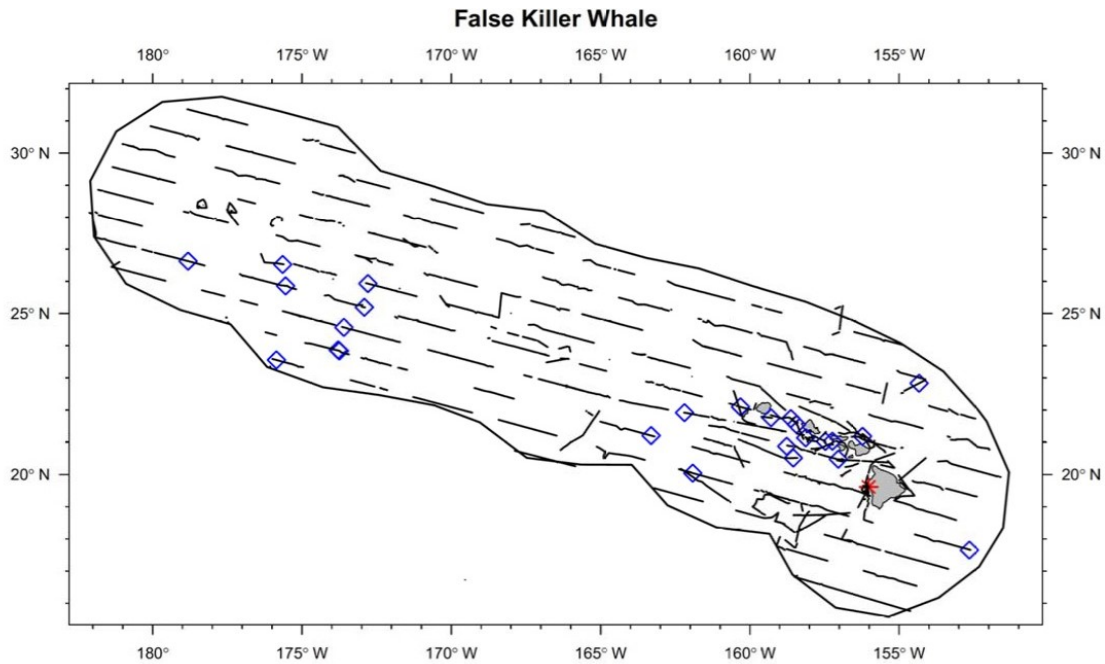


Figure 1.2: A map by [Yano et al. \(2018\)](#) showing the survey area and False Killer Whale sightings. The survey includes the waters immediately surrounding the Hawaiian Islands in the lower right, but also extends further out to the PMNM area to the upper left. The black diagonal lines show the visual effort which coincide with the transects of the ship journeys. The red asterisks show visual detection only of False Killer Whales, while the blue diamonds indicate both acoustic and visual detections.

The survey involved the two research ships *Oscar Elton Sette* and *Reuben Lasker* independently travelling along transects (equidistant, fixed and straight paths) in the survey area for four weeks at a time, with 46 scientists assigned between the two ships. The survey effort involved both visual and audio assessments of *cetaceans* (whales and dolphins) and seabirds. Aside from the visual and audio survey, biopsy samples were collected and satellite tags deployed on cetaceans.

During the survey, three groups of False Killer Whales were observed: *Insular*, which inhabit the waters between and near the main Hawaiian islands; *Northwestern*, which live mostly in the PMNM area; and *Pelagic*, which live in the open ocean waters near the main Hawaiian islands and the PMNM ([Carretta et al. \(2022\)](#)). While all three groups have been found to be demographically distinct through genetic, telemetry, and photo-identification studies, their geographic ranges overlap each other around the Kauai and Niihau islands, which are the Northwestern-most main Hawaiian islands.

This study focuses on *Insular* and *Pelagic* groups only. As of 2015, the abundance estimate for the Hawaiian *Insular* stock was 149 animals, while the *Pelagic* stock was 928 animals (Carretta et al. (2022)). In other words, there are substantially fewer *Insular* animals than *Pelagic*, and out of 1,077 animal counts, it is expected that 13.8% would be *Insular*, while 86.2% would be *Pelagic*.

### 1.3 Audio Data Collection and Curation

The original audio data for this study was taken from HICEAS 2017 was created by recording audio using towed hydrophone arrays while the ships were running along transects during the survey (Yano et al. (2018)). These files were then curated by the University of St Andrews - Dolphin Acoustics Vertically Integrated Project (VIP) team (McCloskey (2024)).

Only audio files that had been positively identified as having *Pseudorca* sounds with no other animal species present were included. The audio files were then split into 1-second clips for the VIP students to trace the whistles using PAMGuard and ROCCA, which are software used for passive acoustic data processing and species identification.

The dataset provided comprised 3 main folders: *Pelagic* with 2,842 audio files; *Insular* with 439 files, and *Unknown* with 4,637 files. This last category contained audio files that came from either group, but could not be positively confirmed. No audio files were provided for the *Northwestern* group.

Comparing the two known groups, audio data from the *Pelagic* group comprised 86.7% of all the tagged data, with *Insular* only accounting for 13.4% of the dataset. This imbalance is to be expected since *Pelagic* False Killer Whales comprise 86.2% of animal counts between the two groups, as discussed in Section 1.2.

Subfolders in each group were descriptively named after the ship or leg journey (*Sette* or *Lasker*), and the audio file were named after the date and time of the recording.

The audio files contained *Pseudorca* whistles and clicks, but also captured ambient audio such as boat engine noise and ocean wave sounds. The files varied from one channel to six audio channels, and ranged from around 0.9Mb to 3.0Mb in size (with a few files larger than 7.0Mb), totalling around 16 gigabytes.

For this study, only the identified groups *Pelagic* and *Insular* were used during model training, while the *Unknown* set was used for the prediction simulation.

## 1.4 Audio Visualisation Using Spectrograms

Due to the large size of the audio files, it was necessary to convert the data from audio to some form of sound visualisation to make the data more manageable.

Sound is perceived when a mechanical action moves through physical matter as vibrations ([Bradley & Stern \(2008\)](#)), and can be described using several terms: *Wavelength* describes the distance of each *wave* (or vibration) from each other; *Frequency* (f) is the number of repetitions (or vibrations) that occur per period (usually one second) and is described with *Hertz* (Hz); and *Amplitude* describes the magnitude or strength of the vibrations.

High amplitude does not necessarily mean high frequency in the audio, but refers rather to the strength of the individual waves. A high frequency file and low frequency file could therefore have the same amplitude, as shown in Figure 1.3.



## FREQUENCY, WAVELENGTH, AND AMPLITUDE

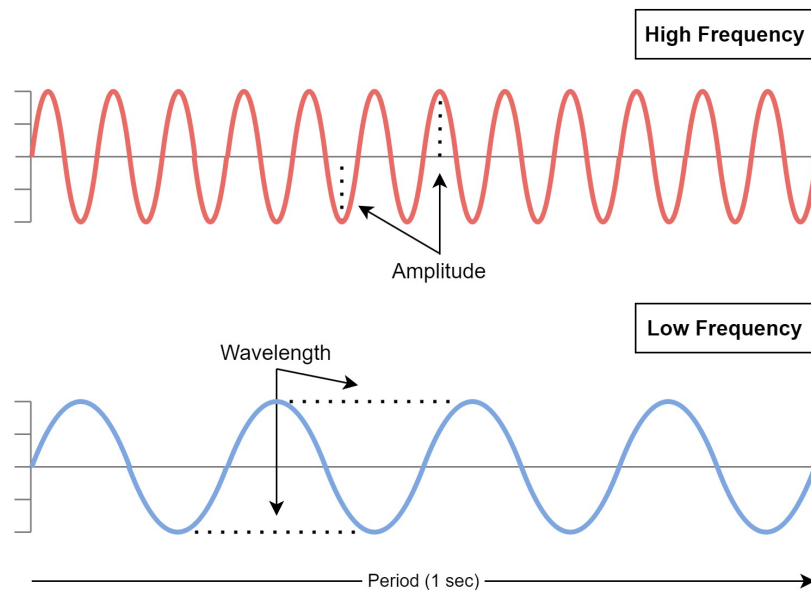


Figure 1.3: A simple diagram showing the relationship between frequency, wavelength and amplitude.

There are several ways to visualise sound; a straightforward way would be to plot the combined amplitude of its frequencies over time as a *Time domain graph* ([Abbot \(2024\)](#)) as in Figure 1.4 (a). Although this shows the strength of the sound over time, it leaves out the intricate details of the sound's frequencies.

Another way would be to show the frequencies in a *Frequency domain graph* as in Figure 1.4 (b). This graph shows the sound's frequencies for a certain period, and it is possible to have several different frequencies on one graph. However, this graph is not able to show the time period associated with the sound.

A different form of visualisation is therefore needed to show all three dimensions of frequency, amplitude, and time.

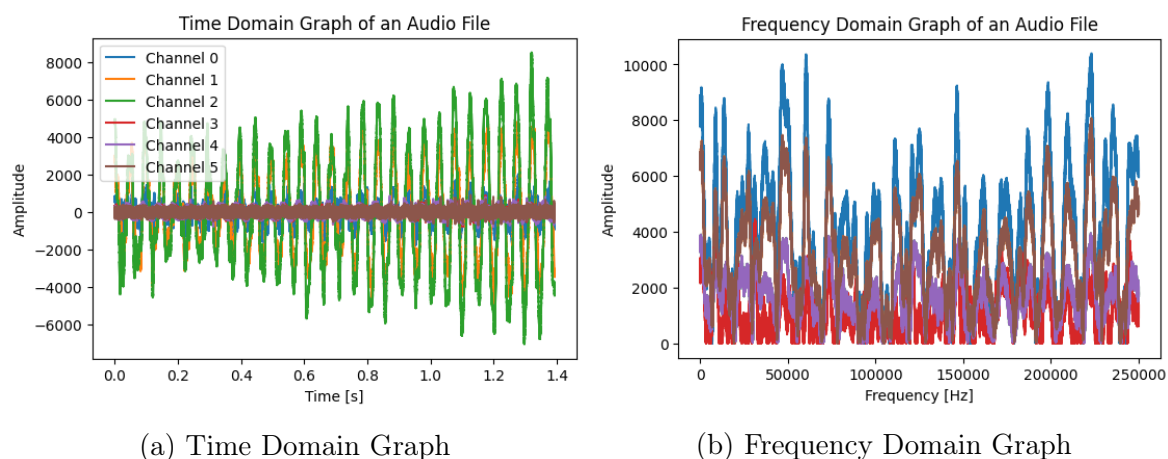


Figure 1.4: An audio file from the *Pelagic* class represented in a time domain (*left*) and the same file as a frequency domain (*right*). The left graph shows the audio amplitude of the signal across time, while the right graph shows the amplitude for the different audio frequencies. The different colours represent the audio channels in the file.

A Spectrogram is a combination of a time domain graph and a frequency domain graph in that it shows frequencies present and their strength over the course of time. Spectrograms are created using a process involving the *Short-Time Fourier Transform* (STFT). This process takes in the amplitudes from the time domain and cuts the data into equal segments across the time period. It then computes the Discrete Fourier Transform per segment to show the magnitude of the different frequencies for that segment, and applies a *window* function to smooth the edges of each segment in the entire spectrogram (Oppenheim et al. (1998)).

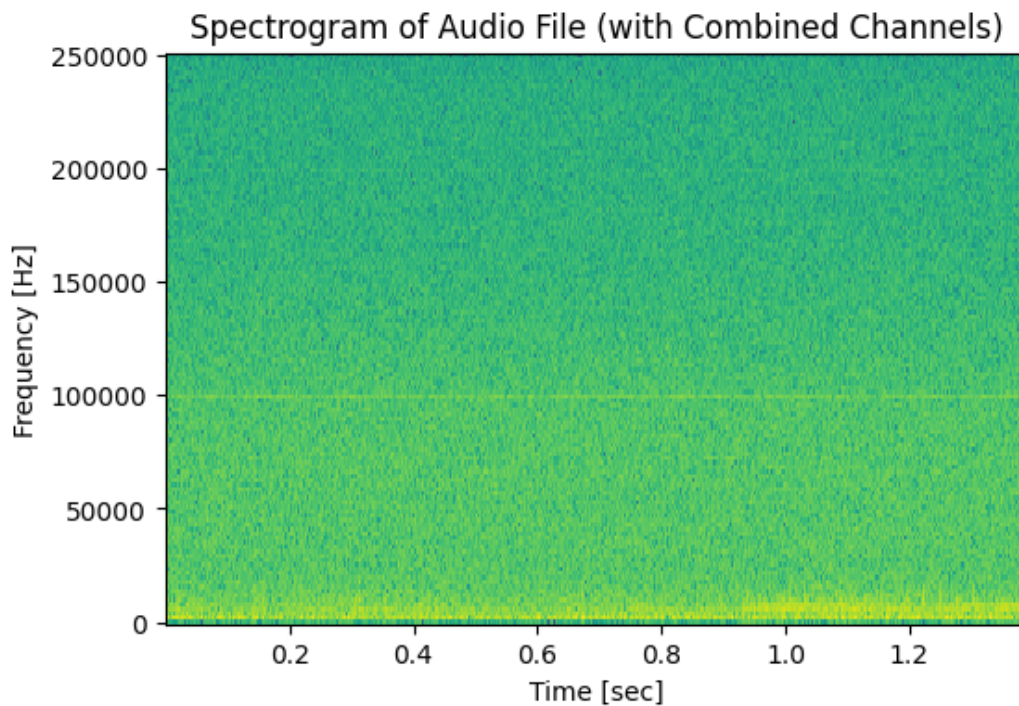


Figure 1.5: The audio file in Figure 1.4 now shown in a spectrogram. Unlike the previous plots, this plot shows the audio frequencies, amplitude and time in one graphical representation. While False Killer Whale whistles should be seen at around 5,000Hz and clicks at around 40,000Hz, this spectrogram shows high amplitude at 10,000Hz throughout the period, as well as around 200Hz starting around 1.0 sec. This might be indicative of other sounds such as engine and wave noises picked up by the hydrophone recorders.

A spectrogram can then be read as a series of frequencies on the y-axis with the amplitude for the different frequencies represented by a colour scale, shown over time on the x-axis.

## 1.5 Motivation for Study

The HICEAS survey is part of NOAA's mandate to monitor the population and habitats of marine mammals around Hawaii ([NOAA Fisheries \(2024\)](#)). These surveys are extremely time-intensive (up to five months continually at sea) and labour-intensive (at least 40 science officers are required excluding the ships' crew) ([Yano et al. \(2018\)](#)).

Accurate identification of the False Killer Whales around the Hawaiian waters is important to properly monitor the state of the groups related to the HICEAS mandate. This information in turn would be of interest to the United States Navy which conducts naval training in the area, as well as the US Bureau of Ocean Energy Management which assesses future sites for wind-warm development ([Yano et al. \(2018\)](#)).

If the False Killer Whale groups can be classified using audio files only, then labour-intensive recording methods could substantially be reduced, and instead other more cost-effective monitoring tools could be used such as bottom-mounted autonomous recorders, which can be deployed for years at a time ([Sousa-Lima \(2013\)](#)), and autonomous gliders, which are extremely energy-efficient and require fewer deployment and recovery intervals ([National Oceanography Centre \(2024\)](#)). Therefore, future data collection could be less time-consuming, less expensive, and less labour-intensive.

## 2 | Literature Review

### 2.1 Machine Learning for Acoustic Data

Prior to machine learning techniques, a variety of manual and computer-aided methods were used for identifying and classifying acoustic data. These include the unaided listening and identifying of audio recordings ([Gerhard \(2003\)](#)) and inspecting spectrograms or other visualisations ([Digby et al. \(2013\)](#)). However, these methods may introduce some human bias by the observer, or may not be feasible for very large datasets ([Honig & Schackwitz \(2023\)](#)). These shortcomings are addressed with machine learning.

Machine learning for audio data typically start with some form of feature extraction. This involves transforming audio data into some vector data which still represent the audio, but in a smaller format ([Giannakopoulos & Pikrakis \(2014\)](#)). Data is usually converted into Spectrograms, Linear predictive coding (LPCs), or Mel-Feature Cepstral Coefficients (MFCCs); LPCs use a linear function on the audio signal ([Bradbury \(2000\)](#)) while MFCCs make use of the Fourier Transform, but with log transform on the mel (or melody-based) scale ([Bäckström et al. \(2022\)](#)). Both LPCs and MFCCs are most used for applications involving human speech because they describe features approximating the human voice, but they also have wider applications in bioacoustics and acoustics in general.

Several studies have successfully made use of MFCCs and different statistical methods for general audio classification. [Ma et al. \(2006\)](#) showed that a Hidden Markov Model (HMM) trained on MFCCs was able to accurately classify ambient office and street sounds with an accuracy of 93% to 96%. A Gaussian Mixture Model (GMM) trained on MFCCs worked well on music classification (also known as Music Information Retrieval) but plateaus for complex variations. [Aucouturier & Pachet \(2004\)](#) showed the R-precision (i.e., the precision measure for information retrieval) struggled to get past 65%.

In bioacoustics, [Cheng et al. \(2012\)](#) used support vector machines (SVMs) on LPCs and MFCCs, and achieved up to 100% accuracy on bird songs from three different species. An HMM used by [Clemins et al. \(2005\)](#) for the identification of African elephant calls showed overall accuracy of 79.7% (but still achieved 90.9% for specific animal vocalisations of rumbles). [Fagerlund \(2007\)](#) used a multi-class classifier using two different datasets of at least six bird species; the classifier involved a decision tree with binary SVM classifiers at each node using MFCCs, and achieved accuracy scores from 91% to 98% depending on the species.

In a review of audio classification techniques used from 2000-2021, [Mutanu et al. \(2022\)](#) showed that the most popular pre-processing techniques include Short-Time Fourier Transformations, MFCCs, and LPCs. Machine learning ensemble algorithms used here resulted in accuracy scores of up to 87%, while Convolutional Neural Networks (CNNs) similarly resulted in accuracy scores of up to 88% (CNNs in bioacoustics are discussed in detail in section 2.3).

While earlier statistical learning techniques have shown good results in bioacoustic classification, later studies show more interest in Deep Learning and Neural Network techniques.

## 2.2 Convolutional Neural Networks

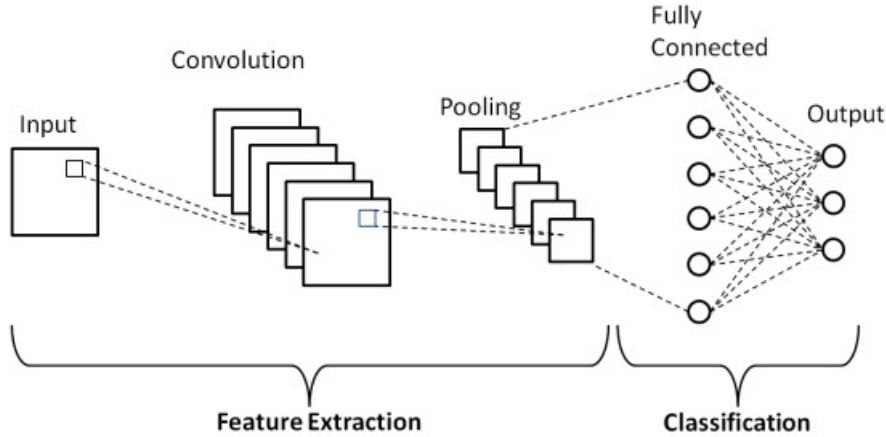


Figure 2.1: A diagram of a simple Convolutional Neural Network by [Phung & Rhee \(2019\)](#). The *Convolution* layer extracts features, while the *Pooling* layer downsamples the data.

Convolutional Neural Networks (CNNs) are a type of neural network that are highly suited for image data as they are able to extract features while reducing the dimensionality of the inputs ([O'Shea & Nash \(2015\)](#)). While CNNs share some properties with other neural networks such as having interconnected node structures, activation functions, and learning through backpropagation, they also have unique layers called *Convolutional* layers and *Pooling* layers ([Géron \(2023\)](#)). Figure 2.1 show these layers in relation to the rest of the neural network.

The *Convolutional* layer extracts important attributes from the input image before passing it on to the next layers. The image input can be thought of as a matrix of pixel data, and the convolutional filter (also known as a kernel) transforms every pixel in all rows and columns of the matrix.

$$I_K(u, v) = \sum_{h=-m/2}^{m/2} \sum_{k=-m/2}^{m/2} K(h, k) I(u + h, v + k)$$

The formula above shows the convolution linear expression described by [Fusiello \(2024\)](#), where  $I_K$  is the output of a convolution kernel  $K$  on an image input  $I$ , with  $(u, v)$  denoting the pixel position of the input and output image, and  $(h, k)$  denoting the row, column position in the kernel, and  $m$  denoting the height or width of the kernel. Figure 2.2 shows how this is applied in a matrix.

1	2	3	4
5	6	7	8
4	3	2	1
8	7	6	5

 $\times$ 

1	2
3	4

 $=$ 

44	64
62	42

Figure 2.2: An application of the convolution linear expression showing a convolution kernel of size 2 x 2 applied on a matrix of size 4 x 4 with a step size of 2-cells both horizontally and vertically. The kernel  $K$  is applied to each colour quadrant of the image  $I$ , and outputs one cell value in  $I_K$ . So  $I_K(1, 1) = (1 * 1) + (2 * 2) + (5 * 3) + (6 * 4)$  which is 44, and the process is applied to the rest of the input matrix.

One issue with applying the convolution kernel in the above example is the output matrix size is smaller than the input due to how the kernel moves through the image. The outer cells are not able to contribute as much as the inner cells, so information is lost in the outer perimeter. *Padding* addresses this problem by adding rows and columns to the edges, thereby increasing the size of the input matrix ([Zhang et al. \(2023\)](#)). Figure 2.3 shows how padding can output the same sized matrix as the input.

The *Pooling* layer downsamples the image matrix which helps to reduce the dimensions of the data while still keeping the important features, effectively outputting a “summary” of the data. Figure 2.4 shows an example of the pooling layer.



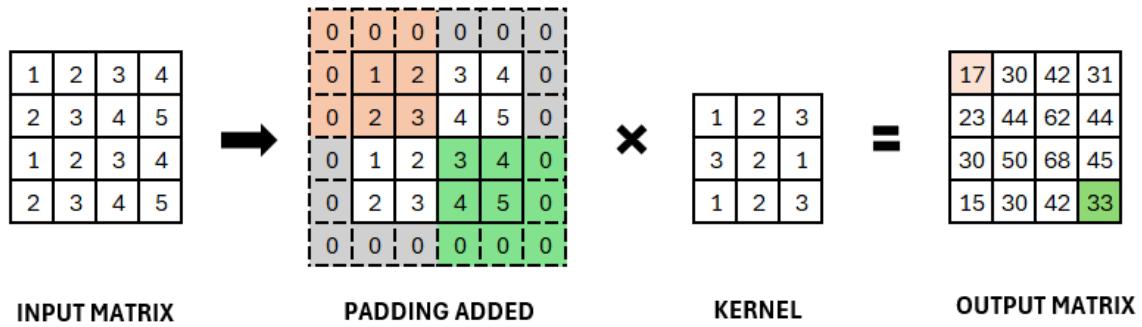


Figure 2.3: A matrix of size 4 x 4 is padded with extra cells around the boundaries. The recommended size of the horizontal and vertical padding  $p_h$  and  $p_w$  with respect to the kernel size  $k$  are  $p_h = k_h - 1$  and  $p_w = k_w - 1$ . Here,  $p_h$  is  $3 - 1 = 2$ , which is the same for  $p_w$ . The padded cells (*in grey*) are distributed around the perimeter and are filled with 0 for all cells (also known as *zero padding*). Here, the step size is 1 cell at a time, so the output matrix will have the same size as the input matrix. Similar to Figure 2.2,  $I_K(1, 1)$  can be computed as  $(0 * 1) + (0 * 2) + (0 * 3) + (0 * 3) + (1 * 2) + (2 * 1) + (1 * 0) + (2 * 2) + (3 * 3) = 17$ , and is similar for the rest of the cells for the output image  $I_K$ .

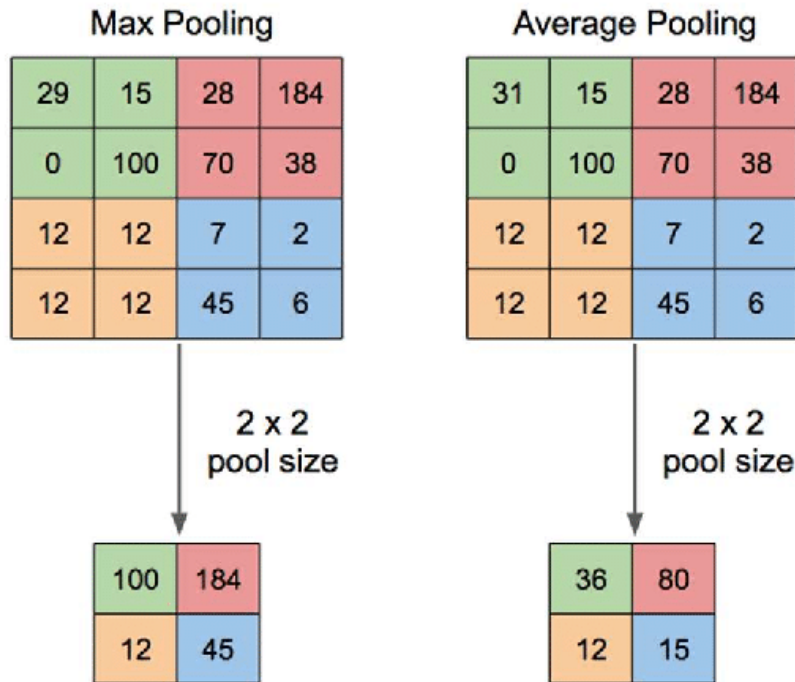


Figure 2.4: Illustrative examples of different pooling layers by [Yani et al. \(2019\)](#), showing how the pooling operation downsamples through the feature map. As their names suggest, Max Pooling (*left*) outputs the maximum value from each patch, while Average Pooling (*right*) outputs the average value for each patch.

One other important concept in convolutional layers is the *Stride*. This dictates how far the kernel steps through the image matrix. In Figure 2.3, the stride value was 1, and with the added padding, the matrix size is maintained. In the earlier example Figure 2.2 however, the stride value is 2, and outputs a smaller sized matrix. The stride can therefore allow the convolutional layer to extract features and downsample at the same time. Most CNN architectures have alternating convolutional and pooling layers, but [Springenberg et al. \(2015\)](#) showed that it is possible to replace the pooling layers with an increased stride value without loss in the CNN’s accuracy.

CNNs will have different architectures for their convolutional and pooling layers, but their basic application is the same; they take inputs and then extract and downsample the inputs before passing them on downstream. Figure 2.5 shows a simplified analogy of the convolution and pooling layers’ application.

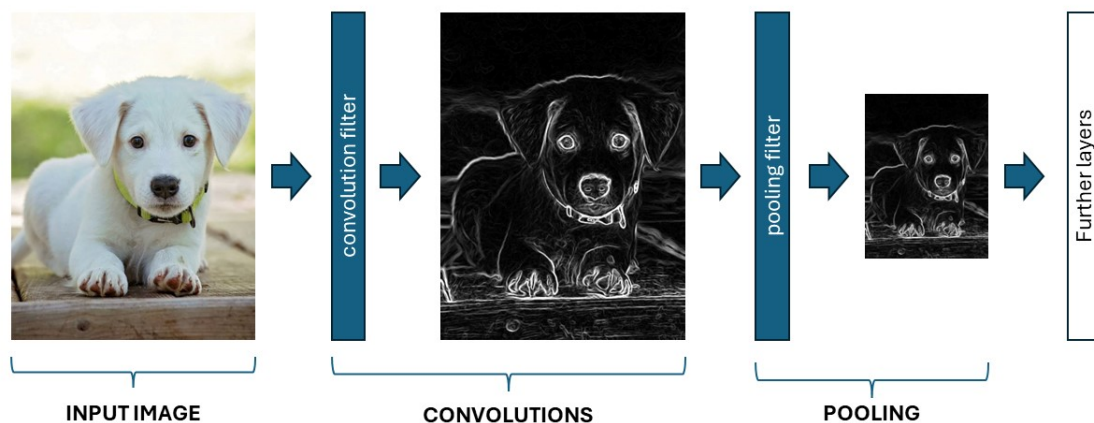


Figure 2.5: An illustrative example of how a convolutional layer and pooling layer process an image of a puppy from [Pixabay \(2016\)](#). The convolutional filter identifies important features such as the puppy’s eyes, nose and paws, and then downsamples the input while still keeping the important features.

While the earliest neural networks were created to address image-based problems such as [Lecun et al. \(1998\)](#)'s LeNet-5 which was used to classify handwritten numbers, some models were applied to the field of bioacoustics. The model by [Parsons \(2001\)](#) was able to identify bat species through echolocation recordings converted into power spectra data. Parson's model was a neural network created using MatLab's NeuralNetwork Toolbox, and had an accuracy rate of 99% for one species.

These models were purposefully built from scratch, but the current practice makes use of pretrained CNNs applied to new and unseen data.

## 2.3 Transfer Learning

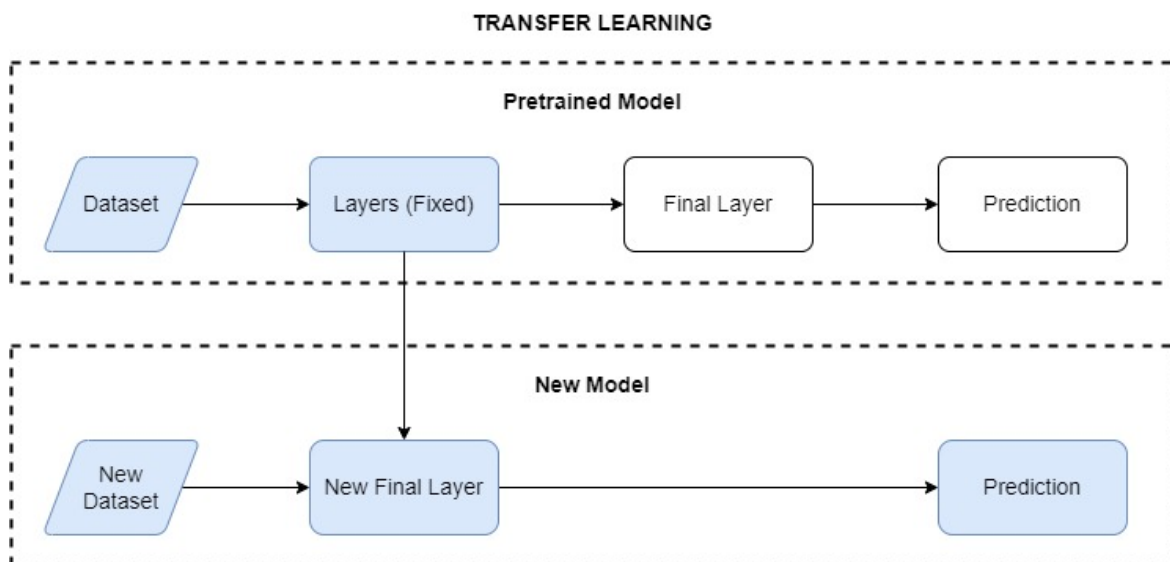


Figure 2.6: A simplified diagram of transfer learning showing how a pretrained model can be used on a different dataset for a new prediction.

Transfer Learning refers to the practice of using a model pretrained on a domain set and using the model as a feature extractor on a different target set ([Chilamkurthy \(2024\)](#)) or by fine-tuning the layers of the model through backpropagation ([Li & Adeli \(2024\)](#)). This process makes use of the “weights” learned by the model on the original dataset, but allows the model to learn on a new dataset, leading to a different prediction as illustrated in Figure 2.6.

Transfer learning also allows the use of smaller datasets which is immensely practical as CNNs typically use large datasets. In this way, researchers do not need to create a model from scratch nor build a large dataset, but are still able to achieve strong results ([Razavian et al. \(2014\)](#)).

Many models are available for transfer learning and have performed well on different datasets than the ones they were originally trained on. GoogLeNet (also known as *Inception*) is a CNN which uses multiple convolution kernels (instead of just one kernel) between layers for increased network depth and width ([Szegedy et al. \(2014\)](#)). Whereas traditional CNNs have connections only to the next layer, DenseNet by [Huang et al. \(2018\)](#) has connections to every other layer, for deeper training. In contrast, ResNet by [He et al. \(2015\)](#) has “shortcut connections” which skip one or more layers while still having sequentially connected layers. EfficientNet by [Tan & Le \(2020\)](#) builds on Resnet, but with additional scaling.

While CNNs and Transfer Learning are usually applied in computer vision use cases, these have also been effectively used in acoustics. The work by [Palanisamy et al. \(2020\)](#) on music genre classification using Densenet showed a result of 92.89% accuracy while [Song et al. \(2017\)](#) used Recurrent Neural Networks (RNNs) achieved 95.8% (RNNs are a related type of neural network that works well on sequential data ([Keren & Schuller \(2016\)](#))). Both studies used the GTZAN music dataset ([Sturm \(2014\)](#)) which contain 1,000 audio files of 30-second clips representing 10 music genres.

In the field of bioacoustic classification, strong results have also been found also using animal sounds. In a study by [Zhong et al. \(2020\)](#) on tropical birds and amphibians audio using ResNet50 originally trained on the ILSVRC-2015 dataset, also known as *ImageNet* ([Russakovsky et al. \(2015\)](#)), the researchers achieved an amazing 99.5% Area Under the Curve by adding a custom loss function and pseudo-labeling. A study by [Baptista & Antunes \(2021\)](#) which also used ResNet and a similar bird and frog dataset scored 91% accuracy. A series of experiments by [Ghani et al. \(2023\)](#) included AudioMAE, a transformer-type model originally trained on AudioSet, a database of generic sounds extracted from Youtube ([Gemmeke et al. \(2017\)](#)). The model was used

on Black-tailed Godwit bird calls and achieved an AUC of over 90.0%. However, it should also be noted that in the same series, the models originally trained on bird and mammal sounds achieved far superior results.

For marine bioacoustics specifically, transfer learning has also been applied with similar success rates. [Nur Korkmaz et al. \(2023\)](#) used VGG16 ([Simonyan & Zisserman \(2015\)](#)) to detect the presence of dolphin whistles with an accuracy of 92.3%. In another study, [Williams et al. \(2024\)](#) created called ReefSet, a large annotated dataset on coral reef sounds. They used BirdNet ([Kahl et al. \(2021\)](#)) on ReefSet and obtained 90.8% AUC. Similar to [Ghani et al. \(2023\)](#)'s experiments, they also tried YAMNet (based on MobileNet ([Howard et al. \(2017\)](#))) and VGGish (based on VGG) originally trained on general audio, but these resulted in much poorer scores. [Lu et al. \(2021\)](#) used AlexNet [Krizhevsky et al. \(2017\)](#) on the Watkins Marine Mammal Sound Database ([Sayigh et al. \(2017\)](#)) and achieved 97.4% accuracy on classification tasks. These studies and their corresponding scores are summarised in Table 2.1.

Author and year	Base model	Target dataset	Score
<a href="#">Williams et al. (2024)</a>	BirdNET	ReefSet	90.8% AUC
<a href="#">Ghani et al. (2023)</a>	AudioMAE	Godwit Calls	>90.0% AUC
<a href="#">Nur Korkmaz et al. (2023)</a>	VGG16	Dolphin whistles in Israel	92.3% accuracy
<a href="#">Baptista &amp; Antunes (2021)</a>	ResNet	Puerto Rican birds & frogs	91.0% accuracy
<a href="#">Lu et al. (2021)</a>	AlexNet	Watkins Marine DB	97.4% accuracy
<a href="#">Zhong et al. (2020)</a>	ResNet	Puerto Rican birds & frogs	99.5% AUC

Table 2.1: A selected list of recent studies in CNN Transfer Learning focusing on bioacoustics. While the target datasets focused on bioacoustics (bird calls, dolphin whistles, etc.), the studies all achieved superior results, even with base models that were originally trained on image data or generic audio data.

In all these studies, audio data was first converted into spectrogram or mel-spectrogram images before being applied to the CNNs, however they achieved transfer learning in different ways. [Ghani et al. \(2023\)](#), [Williams et al. \(2024\)](#), and [Lu et al. \(2021\)](#) removed or replaced some of their CNN's layers. [Zhong et al. \(2020\)](#) and [Nur Korkmaz et al. \(2023\)](#) used a custom loss function or optimiser. [Baptista & Antunes \(2021\)](#) experimented with different window sizes in their Mel spectrograms.

Despite the overall good results, the studies by [Ghani et al. \(2023\)](#) and [Williams et al. \(2024\)](#) also highlight that original training done using bioacoustic data (rather than general acoustic sounds or images) still lead to higher scores.

Based on these studies, transfer learning can therefore be achieved by slightly modifying the CNN architecture, or by hyperparameter tuning, or by data augmentation at the source, and respectable scores can still be obtained with CNN models trained on general datasets.

## 2.4 CNN Software Frameworks

Due to the complexity of these models' architectures, software frameworks have been made available for easier research and development in CNNs. Rather than creating these models from scratch, these frameworks allow the use of design patterns (i.e. repeatable “blueprints” for writing code) and CNN models ready for transfer training.

In Deep Learning applications, tensors are used which are a multi-dimensional type of data object ([Kolecki \(2002\)](#)) that allows for faster training on a GPU ([Pytorch \(2024\)](#)). This data structure is used for encoding and decoding model outputs by software frameworks such as TensorFlow and PyTorch.

TensorFlow was developed by Google's internal research team ([Abadi et al. \(2016\)](#)) and released in 2015 as an open-source Deep Learning library, with deep interactions with Google Cloud's platform. PyTorch was developed by Meta AI also as an open-source Deep Learning library but with a more “Pythonic programming style” ([Paszke et al. \(2019\)](#)).

## Frameworks

Paper Implementations grouped by framework

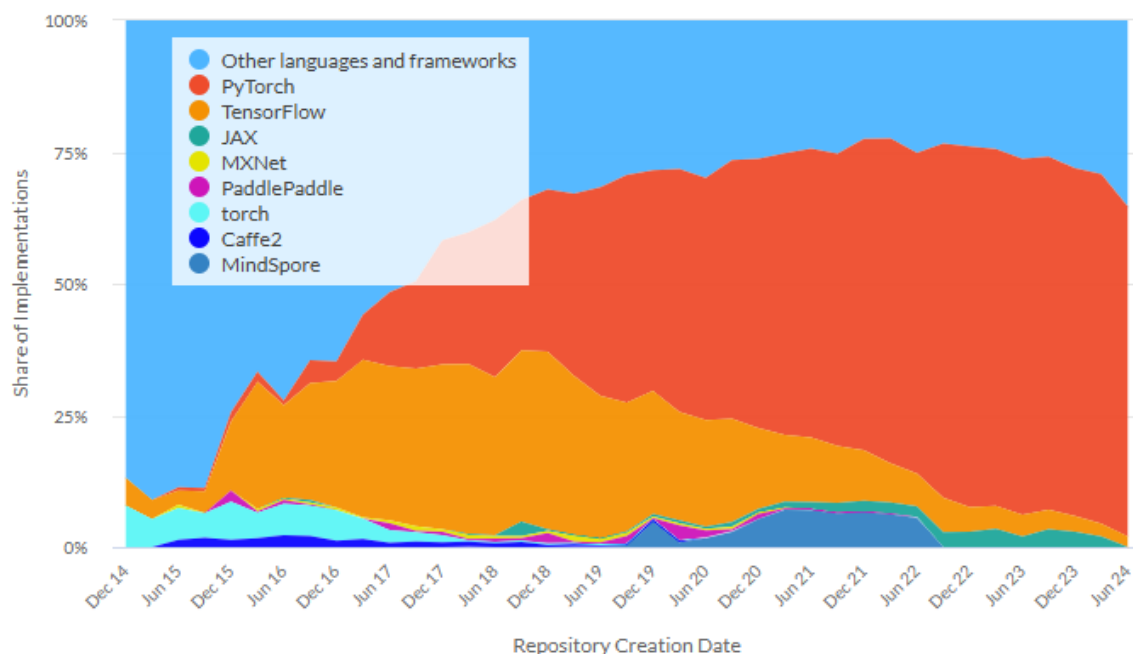


Figure 2.7: A screenshot from *Papers With Code: Trends* (2024) showing the framework usage from June 2014 to June 2024. Here, TensorFlow led the trend up until 2019 when PyTorch took a significant lead.

The graph in Figure 2.7 shows that as of June 2024, *Papers With Code: Trends* (2024) recorded that 63% of published papers' repositories used PyTorch, while only 2% used TensorFlow (and 35% used other languages). PyTorch's popularity along with its "pythonesque" syntax make it ideal for this study's experimentation in CNNs.

## 3 | Methodology

The experiments for this study comprised five main sections; first, the wav files needed to be converted to spectrogram images; second, transfer learning with the image dataset was done using pretrained CNNs; third, the disproportioned dataset was addressed; fourth, the model hyperparameters were manually tuned; and fifth, the best model was used to make predictions on the *Unknown* data. Figure 3.1 shows the flow of these experiments.

This study used the Python programming language for data processing, model training and analysis. Training was done using a remote NVIDIA A30 Tensor Core GPU, configured to 4 cores and 16GB of memory, using a Python virtual environment.

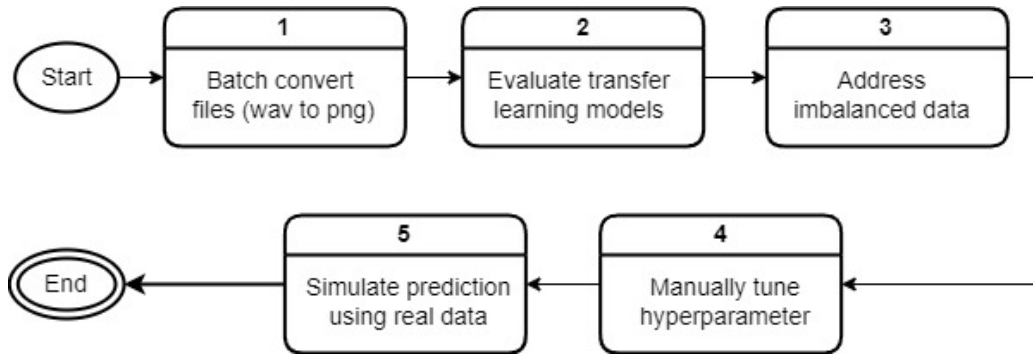


Figure 3.1: Outline of the processes done for this study.

### 3.1 Converting Audio WAVs to Spectrogram PNGs

As previous studies have shown that CNNs work well on image data converted from audio, this study will make use of spectrogram images from the original audio dataset.

A batch function was created that walked through all folders and subfolders in the dataset and checked only folders for *Insular* and *Pelagic* parent folders, ignoring the *Unknown* folder. The function also looked for *wav* files only, ignoring other file types.



For each *wav* file found, the raw audio data was extracted and “flattened” to only one channel. The audio sample data and sample rate were then converted to a *NumPy* ([Harris et al. \(2020\)](#)) array. Using these samples and sample rate, the raw frequencies, times, and spectrogram data were then extracted.

Finally, spectrogram plots were created using the frequencies, times and spectrogram data. The plot size was set to 12 x 8 inches for easier visual inspection, and were created without any axes and labels as these are not important to the CNN model. The created plots were then saved to their respective *Insular* and *Pelagic* folders.

## 3.2 Transfer Learning Using Pretrained Models

Due to the small dataset size of around 2,800 audio files, transfer learning was used with pretrained CNNs as fixed feature extractors (i.e., the weights in the network are fixed except for the final fully connected layer only). Transfer learning was done using PyTorch and *Torchvision* ([2024](#)) as their “pythonesque” interface flowed well with the rest of the Python scripts, and the PyTorch code by [Chilamkurthy \(2024\)](#) was readily available for use. Code Listing 3.1 shows the relevant snippet where the pretrained model can be specified and the last layer only can be updated.

```
1 # Loading the pretrained model
2 model_conv = torchvision.models.resnet18(weights='IMAGENET1K_V1')
3
4 # Freezing the model's other layers
5 for param in model_conv.parameters():
6     param.requires_grad = False
7
8 # Changing the number of features to 2 (for Insular and Pelagic)
9 num_ftrs = model_conv.fc.in_features
10 model_conv.fc = nn.Linear(num_ftrs, 2)
11
12 # Optimising the parameters of the final layer "fc" layer only
13 optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=0.001,
                             momentum=0.9)
```

Code Listing 3.1: Code snippets adapted from the official PyTorch documentation by [Chilamkurthy \(2024\)](#) showing how to use a CNN as fixed feature extractor in PyTorch. This example shows ResNet, but the code is similar for the other CNNs in this study.

The pretrained models Densenet, Resnet, Googlenet and EfficientNet were used as base models as these were used in similar studies that showed good results, and they share similar APIs making the code easier to configure between the pretrained models. All the models were loaded from *torchvision.models*, and *IMAGENET1K\_V1* was used for the weights.

The number of features in the final layer was changed to “2” for the target classes (*Pelagic* and *Insular*), and only this fully connected layer was optimised while the rest of the layers were fixed.

The optimiser used was the same as the default provided in the Pytorch documentation (SGD, learning rate set to 0.001, and momentum set to 0.9), and the image inputs transformed to 224x224 pixels (the minimum required for all models except for EfficientNet which required 384x384), converted to tensor, then normalised with mean values of [0.485, 0.456, 0.406], and standard deviation values of [0.229, 0.224, 0.225].

The converted images were then loaded by batch using Torchvision’s Datasets and Dataloader, randomly split into train, validation and test sets of 60%, 20%, 20% respectively (with the test set only used on the best model during prediction simulation). However, although these experiments used the *random\_split* method, the *torch.manual\_seed* was set to “220029955” to ensure reproducibility across all experiments.

The models were trained in 24 epochs using a batch size of 64, and the loss and accuracy were computed for each training epoch, with the highest scoring epoch model saved. The accuracy scores for all models were then compared to determine the best pretrained model to be used for the rest of the study.

Although accuracy might not be ideal for imbalanced data, it was used here to compare with related studies’ accuracy scores as discussed in Section 2.3. To provide a better gauge of performance, the precision and recall metrics were also included as these take into account how the model predicts for the underrepresented class.

### 3.3 Addressing Imbalanced Data

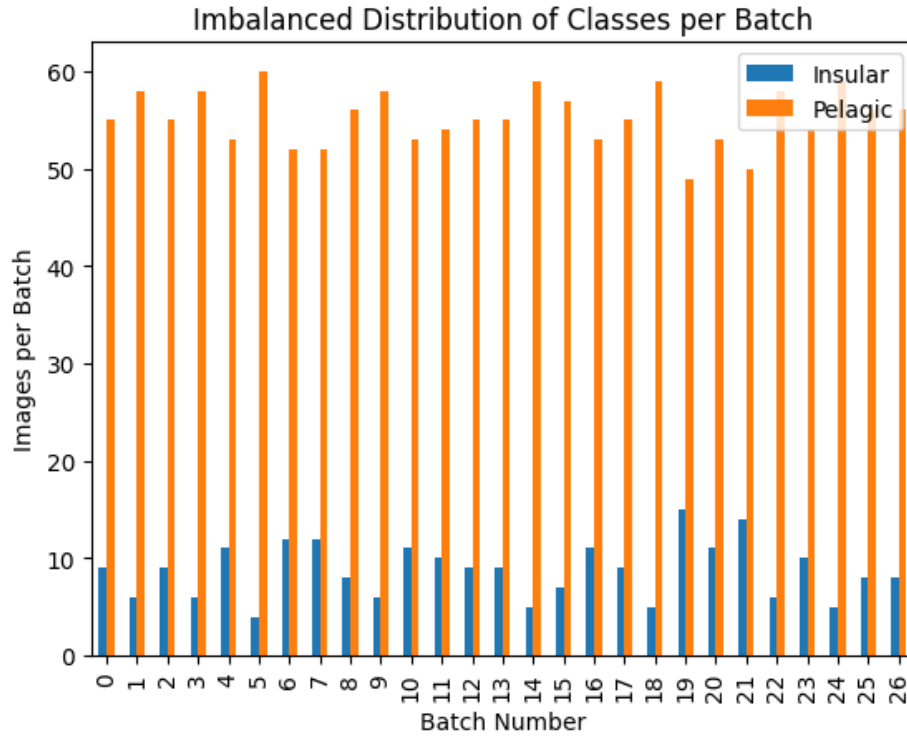


Figure 3.2: Imbalanced data per batch using PyTorch’s Dataloader shows the *Pelagic* has overwhelmingly more selected instances than *Insular*.

The dataset was severely imbalanced with audio from the *Pelagic* group comprising 86.7% of all data. This imbalance led to the Dataloader continuously undersampling the *Insular* class for all batches as shown in Figure 3.2. Predicting without addressing this imbalance would lead to a poorly performing predictor (Yadav & Bhole (2020)), or as Géron (2023) calls it, a *weak learner* (i.e., a classifier that performs as well as random guessing). The CNN could score an accuracy of close to 86.7% which looks good without context, but would be the same result as the model guessing all instances as the majority class *Pelagic*.

An alternative here is to use straightforward undersampling of the majority *Pelagic* class to match the number of *Insular* instances. However, this would lead to severe loss of information as there would only be 439 data points for each class.

To address this, the *Weighted Random Sampler* from [torch.utils.data \(2023\)](#) was used which assigns a weight to each item in the dataset for the Dataloader sampler, allowing the Dataloader to undersample the *Pelagic* class, and oversample (sample with replacement) the *Insular* class. While oversampling may lead to overfitting due to the synthetic instances used ([Alkhawaldeh et al. \(2023\)](#)), it is still a better alternative as using the imbalanced data would lead to the model learning only the majority class, and using straightforward undersampling would lead to loss of information.

In addition to the accuracy, precision and recall metrics, Precision-Recall Curves and Precision-Recall Threshold plots were also used to evaluate the effectiveness of the weighted sampler.

### 3.4 Hyperparameter Tuning and Final Testing

Hyperparameter tuning allows for changing the model values that affects the training process, leading to a better performing model. Instead of external software libraries, hyperparameter tuning for this study was done manually by looping through different values for the *learning rate* and *momentum*.

First, a python script looped through the following learning rates with the momentum fixed at 0.90: [0.0005, 0.001, 0.005, 0.01 0.5, 0.1]. Accuracy scores were compared to determine the best learning rate. Then, with the best learning rate, another script looped through the following momentum values: [0.90, 0.95, 0.97, 0.99].

The accuracy scores are then compared and the hyperparameters with the best scores are selected as the best model.

As a final test, the best model was compared against the unseen test set. A confusion matrix was also created to show how the model performs for the minority and majority classes.

## 3.5 Simulating Predictions with Real Data

```
1 output:  tensor([[ 0.3622, -0.5187]])
2 index: 0
3 class name: Insular
4 probabilities: tensor([[0.7070, 0.2930]])
5
6 There's a 70.7 % chance that this audio file belongs to the Insular
  group.
```

Code Listing 3.2: Sample output when predicting for each test audio file.

To prototype the best model in a realistic setting, a simple Python script was created which takes in a wav file, creates a spectrogram plot, transforms the spectrogram into a tensor, and inputs that tensor into the model.

The model then outputs the predicted class along with the probability for that class.

10 audio files were then selected from the *Unknown* dataset through convenience sampling (i.e., the samples were selected randomly by the researcher).

## 4 | Results

A total of 2,842 spectrogram plot images were created from the audio dataset without any labels or axes marks so only spectrogram image data was fed into the CNN, as shown in Figure 4.1. Each spectrogram file was around 300kb in size, and with dimensions of 950 x 636 pixels (12 x 8 inches) each, although these were resized later during the transfer learning to 224 x 224 or 384 x 384 depending on the pretrained CNN.

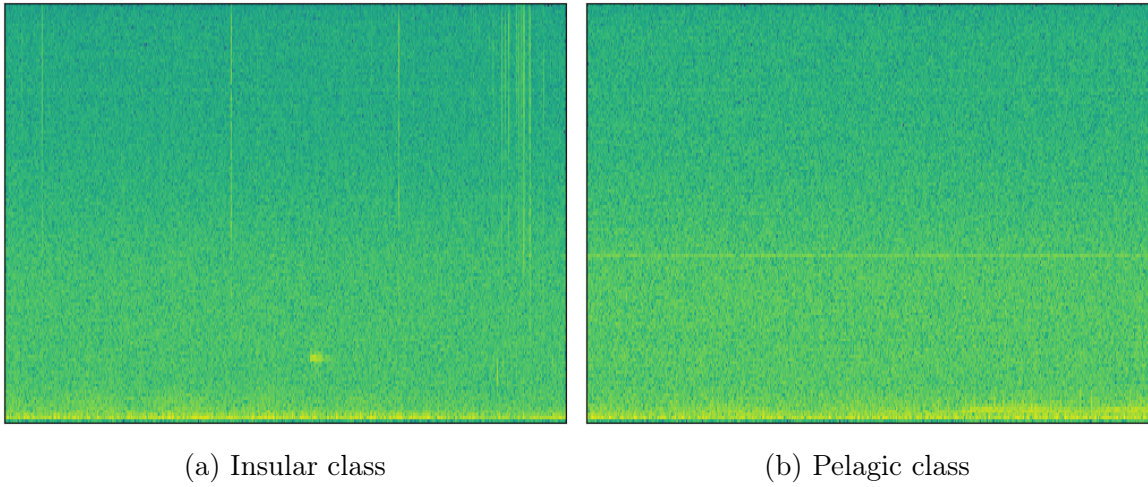


Figure 4.1: Spectrograms from the first wav files in each class; *Insular* on the left and *Pelagic* on the right. The vertical streaks in the *Insular* example are not found in the *Pelagic* class. In contrast, the horizontal line in the middle of the *Pelagic* example can't be seen in the other class. However, these may be attributed to other ambient ocean sounds.

The different pretrained CNN models were loaded and transfer learning applied on the spectrogram imageset, with ResNet showing the best accuracy score of 0.9238. DenseNet followed with an accuracy of 0.9149, but GoogLeNet and EfficientNet performed poorly with accuracies below 0.9, and both had recall values of 1.0. A summary of their results can be seen in Table 4.1.

Base Model	Accuracy	Precision	Recall
ResNet	0.9238	0.9273	0.9918
DenseNet	0.9149	0.9122	0.9980
GoogLeNet	0.8759	0.87479	1.0000
EfficientNet	0.8670	0.8670	1.0000

Table 4.1: Validation results from transfer training using different pre-trained models, sorted by best accuracy score descending.

These accuracy scores resulted from using the imbalanced data, and even ResNet showed poor results for its Precision-Recall tradeoff as seen in figure 4.2. The precision-recall curve should show the line curving toward the 1.0 precision and 1.0 recall (in the upper right), and the precision-recall threshold should show the precision and recall lines meeting at around 0.5. However, neither was the case for these graphs, indicating a critical need to address the imbalance.

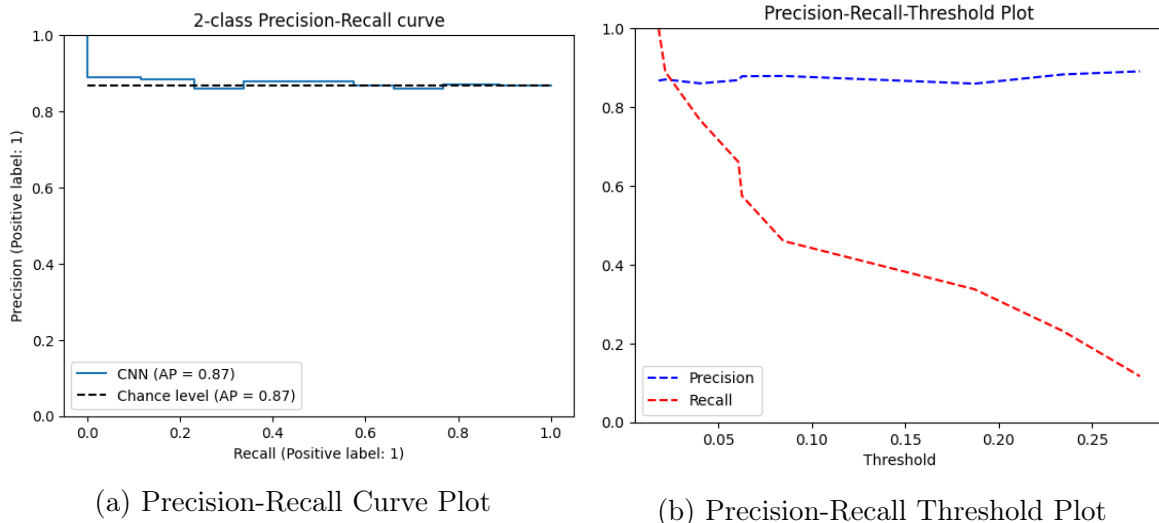


Figure 4.2: The precision-recall curve (*left*) show how the model learns towards a precision of 0.87, which is about the same as the majority class ratio of 87.6%. The threshold plot (*right*) also show the effect of the severely imbalanced data, with the threshold at close to 0.0 where recall is close to 1.0. These plots show the effect of the imbalanced majority class in the dataset.



The ResNet transfer learning model was run again with weighted data as shown in Figure 4.3. With the *Weighted Data Sampler* applied, the ResNet model was run with the balanced Datasampler, and the resulting model yielded a lower accuracy score of 0.9078 compared to the previous imbalanced dataset as seen in Table 4.2.

Model	Accuracy	Precision	Recall
ResNet (Imbalanced)	0.9238	0.9273	0.9918
ResNet (Weighted)	0.9078	0.9700	0.9223

Table 4.2: Validation accuracy results comparing the ResNet model *before* and *after* balancing the dataset. The accuracy dipped slightly after applying Pytorch’s *Weighted Random Sampler*.

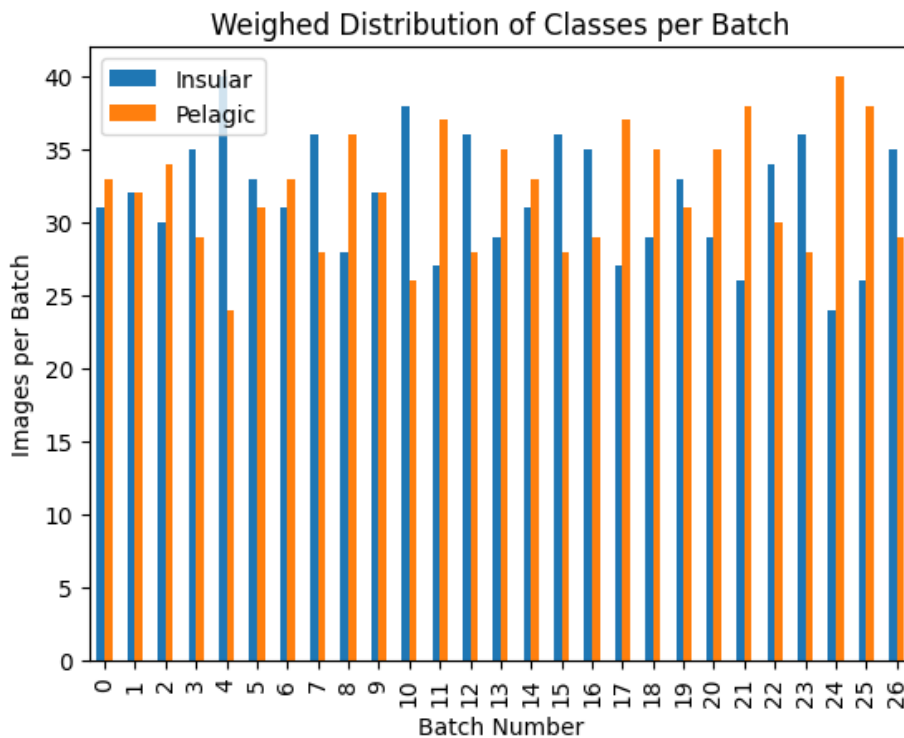


Figure 4.3: Using Pytorch’s *Weighted Data Sampler*, the Dataloader was able to undersample *Pelagic* instances and oversample *Insular* instances which produced a more balanced dataset across all batches.

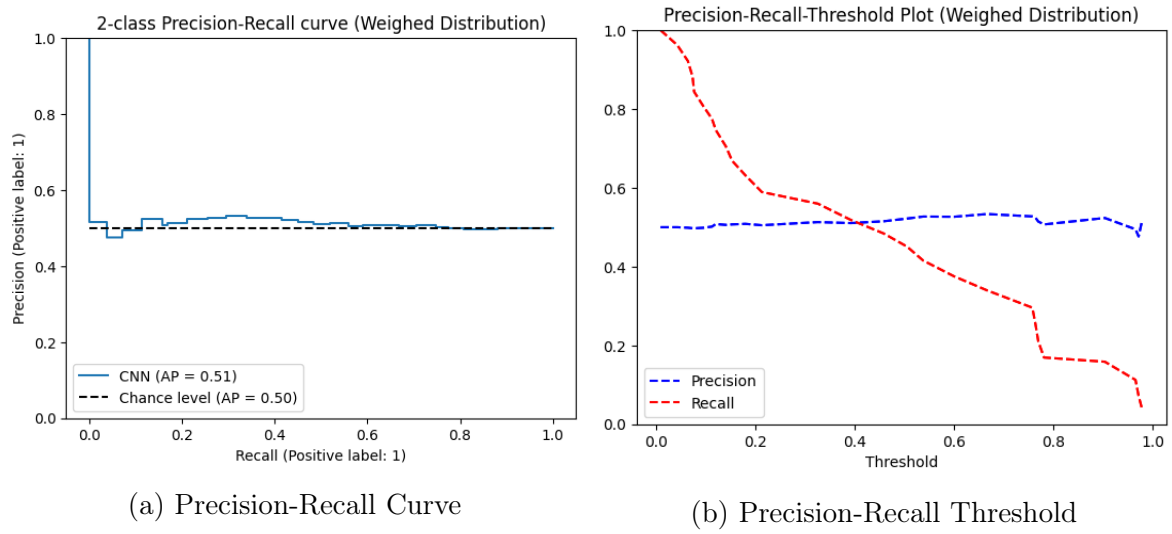


Figure 4.4: ResNet Precision-Recall plots after applying *Weighted Data Sampler*. The precision and recall threshold (*right*) came closer to a more balanced threshold of around 0.4, however the precision-recall curve (*left*) while showing a balance chance level of 0.5 still showed the model line “hugging” the chance line, indicating there were still deficiencies with the model.

Despite the resulting lower accuracy score and poor precision-recall curve, the precision increased significantly to 0.9700 after using the weighted dataset, so hyperparameter tuning was performed on the newer model, with the learning rate tuned first followed by the momentum. The results in Table 4.3 showed the best accuracy scores were a learning rate of 0.01, and momentum values of 0.95 and 0.97.

Loss curve graphs were also created for these, and it was expected that they would show smooth curves for training and validation. However, Figure 4.5 shows erratic movement for both hyperparameters, and was most noticeable for the momentum training.

Learning Rate (LR)		Momentum	
LR Value	Accuracy	Momentum Value	Accuracy
0.0005	0.9131	0.90	0.9415
0.001	0.9291	0.95	0.9486
0.005	0.9450	0.97	0.9486
0.01	0.9486	0.99	0.9468
0.05	0.9433		
0.1	0.9220		

Table 4.3: Accuracy results for manual hyperparameter tuning for Learning Rate (*left*) and Momentum (*right*). For Learning Rate, the momentum was fixed at 0.90, and subsequent Momentum tuning with the best LR of 0.01. The highlighted rows show the best scores.

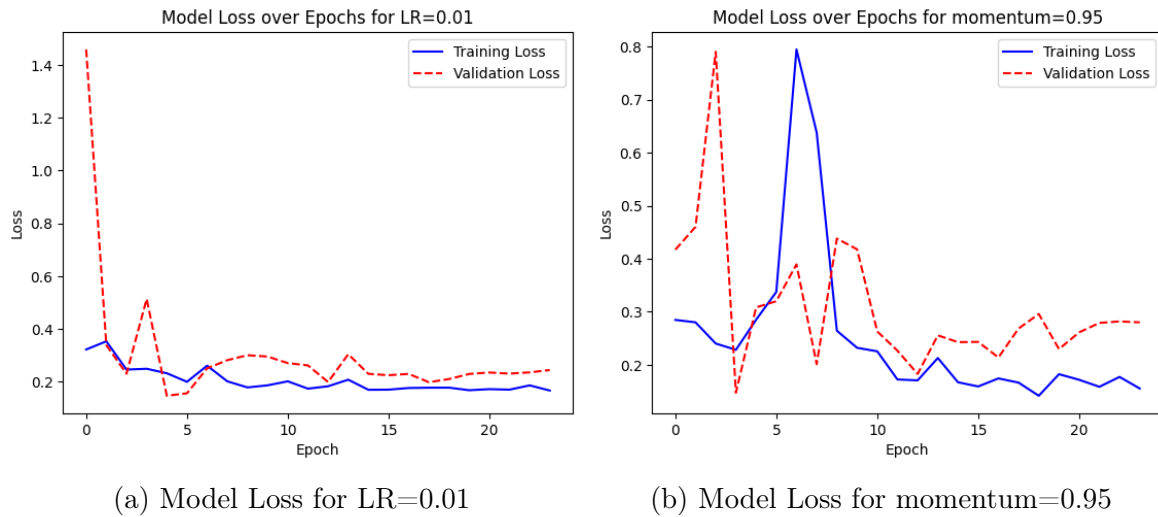


Figure 4.5: Model Losses over Epochs for best learning rate (*left*) and momentum (*right*) over 24 epochs. Despite having the best scores for hyperparameter tuning, the graphs show erratic behaviour.

The final test with the unseen set was done using the ResNet model with the balanced dataset, a learning rate of 0.01 and a momentum of 0.95, resulting in a final accuracy score of 0.9309 as seen in Table 4.4.

(a) Confusion Matrix			(b) Metrics Test Set	
Actual	Predicted		Metric	Score
	1	0		
	1	0		
	48	26	Accuracy	0.9309
	13	477	Precision	0.9483
			Recall	0.9735

Table 4.4: The confusion matrix (*left*) showed the model prediction performance, where the True Positives represent the *Insular* class, while True Negatives are the *Pelagic* class. The model performed a decent job as it correctly identified almost two-thirds of the minority class *Insular*.

Putting these results in the context of the flow of the study, Table 4.5 shows relatively stable accuracy scores from the starting base model’s validation accuracy of 0.9238 up to the final model’s test accuracy of 0.9389.

Table Ref.	Description	Accuracy	Precision	Recall
4.1	ResNet with transfer learning	0.9238	0.9273	0.9918
4.2	Above model + balanced Dataloader	0.9078	0.9700	0.9223
4.3	Above model + best hyperparameters	0.9486	0.9617	0.9755
4.4	Above model on unseen test set	0.9309	0.9483	0.9735

Table 4.5: Final accuracy results comparing the best base model, the base model with weighted dataset, the weighted dataset model with the best hyperparameters, and the same model with the unseen data. Note that all model scores were tested using the validation set, except the final model which used the test set.

The final model was used on a sample of 10 files in the *Unknown* audio dataset, and was able to predict their respective classes as seen in Table 4.6.

Source Folder	Predicted Class	Probability
Lasker_AC_67	Insular	0.5603
Lasker_AC_67	Pelagic	0.9424
Lasker_AC_150	Pelagic	0.8420
Lasker_AC_276	Pelagic	0.9082
Lasker_AC_276	Pelagic	0.9986
Lasker_AC_276	Pelagic	0.9422
Sette_AC_44	Pelagic	0.9770
Sette_AC_91	Pelagic	0.9982
Sette_AC_256	Insular	0.9937
Sette_AC_256	Insular	0.9497

Table 4.6: Predictions using a sample of the *Unknown* audio files. Most of the predictions were given with high probabilities of  $>0.90$ .

## 5 | Discussion

During the initial CNN transfer learning iterations, the model was expected to perform with an accuracy of at least 0.8670 (which is the percentage of *Pelagic* audio in the dataset). ResNet was the best model, and its accuracy of 0.9238 meant that the CNN performed better than a weak learner, and was comparable to the related studies in CNN Transfer Learning (Section 2.3) where one model had an accuracy of 91.0%.

In contrast, GoogLeNet and EfficientNet both had accuracy scores very close to the ratio of 0.87, and had recall of 1.0, suggesting that these models did not perform any better than guessing. Table 5.1 shows how poorly the worst model, EfficientNet, compared with the best model, Resnet. While ResNet had the best results and was selected as the base model for the succeeding experiments, it was interesting to note that the poorly performing EfficientNet was actually built on top of ResNet, and its “compound scaling method” architecture did not work for this use case.

		Predicted	
		1	0
Actual	1	0	75
	0	0	489

(a)

Accuracy: 0.8670  
Recall: 1.0000

		Predicted	
		1	0
Actual	1	37	38
	0	4	485

(b)

Accuracy: 0.9256  
Recall: 0.9918

Table 5.1: Confusion matrices for (a)**EfficientNet** on the left and (b)**ResNet** on the right. Aside from the EfficientNet’s poorer accuracy, it also had a Recall of 1.0 and its confusion matrix shows why this is the case; it predicted all cases as *Pelagic* (True Negatives), while none for *Insular* (True Positives).

The results with the weighted dataset were also interesting since the accuracy decreased after using the *Weighted Random Sampler*, despite the dataset being more balanced. However, the decreased accuracy was not indicative of overall model performance. Table 5.2 below shows how “sacrificing” the majority class *Pelagic* (here, the True Negatives) can lead to more correct cases of the minority class *Insular* (the True Positives), despite the lower overall accuracy.

The increase in precision to 0.9700 also showed that, although the accuracy was lower, overall the model predicted more correct cases for each class. Therefore, the weighted dataset was still used in the next phase of the study.

		Predicted	
		1	0
Actual	1	37	38
	0	4	485

(a)

Accuracy: 0.9256  
Precision: 0.9273

		Predicted	
		1	0
Actual	1	61	14
	0	38	451

(b)

Accuracy: 0.9078  
Precision: 0.9700

Table 5.2: The confusion matrices of the ResNet model (a)**before** and (b)**after** applying the *Weighted Random Sampler (WRS)*, showing how a model with lower accuracy might be more acceptable. In the left matrix, the model has a higher accuracy of 0.9256 with more True Negatives (*Pelagic*), but with fewer True Positives (*Insular*), leading to a lower precision of 0.9273. In the right matrix, the model has a slightly lower accuracy of 0.9078 due to fewer True Negatives (*Pelagic*), but more True Positives (*Insular*), leading to a much higher precision of 0.9700.

The predictions on the *Unknown* dataset shows the feasibility of this model as a prototype binary classifier. Most of the *Unknowns* were found to be *Pelagic*, although this could be due to the weighted data being oversampled for *Insular*, and undersampled for *Pelagic* (i.e., the model might still have learned better with the varied *Pelagic* data).

Based on the above predictions on the *Unknown* set, an assumption could be made that files in the same folder are mostly of the same class, but based on the audio survey methods, one folder may contain recordings from different dates and times, so different groups may still be found in the same folder.

While the final scores for accuracy, precision and recall were overall good, attention should also be given to the precision-recall graphs in and the loss curves in show deficiencies in the model; Figure 4.4 shows the model line still hugging the chance line, and Figure 4.5 shows fluctuations in both training and validation losses. These problems could have been caused by the repetitive synthetic data or some gradient issues,

and could have been addressed with regularisation ([James et al. \(2023\)](#)), optimiser selection ([Géron \(2023\)](#)), or additional hyperparameter tuning.

Despite these shortcomings, this study shows that CNN Transfer Learning can be quite capable of bioacoustics binary classification, and the results suggest that it is possible to distinguish between the *Insular* and *Pelagic* groups of Hawaiian False Killer Whales based on audio data alone.

However, the difference in audio might not be solely attributed to the animal sounds, as other factors might have influenced the CNN classifier. The classifier may have mistakenly considered engine and ambient noise as important features, and sea sounds may have varied due to the effect of water temperature, salinity and recording depth on the recording process ([Bradley & Stern \(2008\)](#)).

On the other hand, if it were to be assumed that the difference could be attributed to the animal sounds alone, then the final test scores can still be improved as similar studies have resulted in much higher metrics (as [Zhong et al. \(2020\)](#) showed an AUC of 99.5%). Therefore, there is still much work that can be done in distinguishing between *Pelagic* and *Insular* False Killer Whale sounds.

Additional methods such as regularisation or additional hyperparameter tuning could have been applied to the CNNs to improve the final models' scores. Based on similar studies, strong results could also be obtained from simpler machine learning methods such as support vector machines which have resulted in 100% accuracy ([Cheng et al. \(2012\)](#)). Alternatively, since the current dataset is small, it could also be possible to attain high accuracy scores by pre-training the CNN on larger bioacoustic datasets such as *ReefSet* rather than on the general dataset *ImageNet* as done by [Ghani et al. \(2023\)](#) and [Williams et al. \(2024\)](#).

Finally, other CNN methods could also be explored with audio and signal processing libraries such as PyTorch's *TorchAudio* ([Hwang et al. \(2023\)](#)) which do not require spectrograms or other image inputs, but instead trains directly on audio data.



## 6 | Conclusion

This study aimed to determine whether it is possible to classify between *Pelagic* and *Insular* groups of Hawaiian Killer Whales based on their audio recordings. The results show that this classification task is possible using simple CNN transfer learning, with an accuracy of 93.89%.

The significant class imbalance in the dataset (with 87.6% belonging to the *Pelagic* class alone) affected the initial transfer learning accuracy scores, but this was addressed using Pytorch libraries. The model’s learning rate and momentum were also tuned to produce slightly better accuracy scores.

The final model was able to predict on an *Unknown* dataset with some confidence, showing that CNN transfer learning can be used effectively for bioacoustic classification.

While this study can be extended with the use of further machine learning and CNN techniques, the current results demonstrate the effectiveness of the CNN “Black Box” with transfer learning; a CNN trained on the *ImageNet* dataset can work remarkably well on *Pseudorcas* audio-to-image dataset by training the last fully connected layer, and achieve strong accuracy, precision and recall scores.

# Bibliography

Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D. G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y. & Zheng, X. (2016), ‘Tensorflow: A system for large-scale machine learning’.

**URL:** <https://arxiv.org/abs/1605.08695>

Abbot, D. (2024), *Understanding Sound*, Pressbooks.

**URL:** <https://pressbooks.pub/sound/>

Alkhawaldeh, I. M., Albalkhi, I. & Naswhan, A. J. (2023), ‘Challenges and limitations of synthetic minority oversampling techniques in machine learning’, *World Journal of Methodology* **13**(5), 373–378.

**URL:** <http://dx.doi.org/10.5662/wjm.v13.i5.373>

Aucouturier, J.-j. & Pachet, F. (2004), ‘Improving timbre similarity: How high’s the sky?’, *J. Negat. Results Speech Audio Sci* **1**.

Bäckström, T., Räsänen, O., Zewoudie, A., Zarazaga, P. P., Koivusalo, L., Das, S., Mellado, E. G., Mansali, M. B., Ramos, D., Kadiri, S. & Alku, P. (2022), *Introduction to Speech Processing*, 2 edn.

**URL:** <https://speechprocessingbook.aalto.fi>

Baird, R. W. (2009), F - false killer whale: *Pseudorca crassidens*, in W. F. Perrin, B. Würsig & J. G. M. Thewissen, eds, ‘Encyclopedia of Marine Mammals (Second Edition)’, Academic Press, pp. 405–406.

**URL:** <https://www.sciencedirect.com/science/article/pii/B9780123735539000973>

Baptista, P. B. & Antunes, C. (2021), ‘Bioacoustic classification framework using transfer learning’, *Model Decision Artificial Intelligence* **35**.

Bradbury, J. (2000), ‘Linear predictive coding’, *Mc G. Hill*.

Bradley, D. L. & Stern, R. (2008), ‘Underwater sound and the marine mammal acoustic

environment: A guide to fundamental principles’.

**URL:** [https://www.mmc.gov/wp-content/uploads/sound\\_bklet.pdf](https://www.mmc.gov/wp-content/uploads/sound_bklet.pdf)

Carretta, J. V., Oleson, E. M., Oleson, E. M., Forney, K. A., Muto, M. M., Weller, D. W., Lang, A. R., Baker, J., Hanson, B., Orr, A. J., Barlow, J., Moore, J. E. & Brownell, R. L. J. (2022), ‘U.s. pacific marine mammal stock assessments: 2021’.

**URL:** <https://repository.library.noaa.gov/view/noaa/44406>

Cheng, J., Xie, B., Lin, C. & Ji, L. (2012), ‘A comparative study in birds: call-type-independent species and individual recognition using four machine-learning methods and two acoustic features’, *Bioacoustics* **21**(2), 157–171.

**URL:** <http://www.tandfonline.com/doi/abs/10.1080/09524622.2012.669664>

Chilamkurthy, S. (2024), ‘Transfer learning for computer vision tutorial - pytorch tutorials 2.3.0+cu121 documentation’.

**URL:** [https://pytorch.org/tutorials/beginner/transfer\\_learning\\_tutorial.html](https://pytorch.org/tutorials/beginner/transfer_learning_tutorial.html)

Clemins, P. J., Johnson, M. T., Leong, K. M. & Savage, A. (2005), ‘Automatic classification and speaker identification of African elephant ( *Loxodonta africana* ) vocalizations’, *The Journal of the Acoustical Society of America* **117**(2), 956–963.

**URL:** <https://pubs.aip.org/jasa/article/117/2/956/541588/Automatic-classification-and-speaker>

Dewynter, M. (2019), ‘Natural science illustrations of marine mammals’.

**URL:** <https://sanctuaire-agoa.fr/editorial/natural-science-illustrations-marine-mammals>

Digby, A., Towsey, M., Bell, B. D. & Teal, P. D. (2013), ‘A practical comparison of manual and autonomous methods for acoustic monitoring’, *Methods in Ecology and Evolution* **4**(7), 675–683.

**URL:** <http://dx.doi.org/10.1111/2041-210X.12060>

DOSITS (2021), ‘False killer whale sounds’.

**URL:** <https://dosits.org/galleries/audio-gallery/marine-mammals/toothed-whales/false-killer-whale/>

- Fagerlund, S. (2007), ‘Bird species recognition using support vector machines’, *EURASIP Journal on Advances in Signal Processing* **2007**(1), 038637.
- Fusiello, A. (2024), *Computer vision: Three-dimensional reconstruction techniques*, Springer Nature, Cham, Switzerland.
- Gemmeke, J. F., Ellis, D. P. W., Freedman, D., Jansen, A., Lawrence, W., Moore, R. C., Plakal, M. & Ritter, M. (2017), Audio set: An ontology and human-labeled dataset for audio events, *in* ‘Proc. IEEE ICASSP 2017’, New Orleans, LA.
- Gerhard, D. (2003), ‘Audio signal classification: History and current techniques’, *Semantic Scholar* .  
**URL:** <https://api.semanticscholar.org/CorpusID:7652753>
- Géron, A. (2023), *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: concepts, tools, and techniques to build intelligent systems*, Data science / machine learning, third edition edn, O’Reilly, Beijing Boston Farnham Sebastopol Tokyo.
- Ghani, B., Denton, T., Kahl, S. & Klinck, H. (2023), ‘Global birdsong embeddings enable superior transfer learning for bioacoustic classification’, *Scientific Reports* **13**(1).  
**URL:** <http://dx.doi.org/10.1038/s41598-023-49989-z>
- Giannakopoulos, T. & Pikrakis, A. (2014), *Introduction to Audio Analysis*, Elsevier.  
**URL:** <https://linkinghub.elsevier.com/retrieve/pii/C20120035247>
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., Fernández del Río, J., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C. & Oliphant, T. E. (2020), ‘Array programming with NumPy’, *Nature* **585**, 357-362.
- He, K., Zhang, X., Ren, S. & Sun, J. (2015), ‘Deep residual learning for image recog-

niton', *CoRR* **abs/1512.03385**.

**URL:** <http://arxiv.org/abs/1512.03385>

Honig, M. & Schackwitz, W. (2023), 'Manual versus semiautomated bioacoustic analysis methods of multiple vocalizations in tricolored blackbird colonies', *Journal of Fish and Wildlife Management* **14**(1), 225-238.

**URL:** <http://dx.doi.org/10.3996/JFWM-22-065>

Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M. & Adam, H. (2017), 'Mobilenets: Efficient convolutional neural networks for mobile vision applications'.

**URL:** <https://arxiv.org/abs/1704.04861>

Huang, G., Liu, Z., van der Maaten, L. & Weinberger, K. Q. (2018), 'Densely Connected Convolutional Networks'. arXiv:1608.06993 [cs].

**URL:** <http://arxiv.org/abs/1608.06993>

Hwang, J., Hira, M., Chen, C., Zhang, X., Ni, Z., Sun, G., Ma, P., Huang, R., Pratap, V., Zhang, Y., Kumar, A., Yu, C.-Y., Zhu, C., Liu, C., Kahn, J., Ravanelli, M., Sun, P., Watanabe, S., Shi, Y., Tao, Y., Scheibler, R., Cornell, S., Kim, S. & Petridis, S. (2023), 'Torchaudio 2.1: Advancing speech recognition, self-supervised learning, and audio processing components for pytorch'.

James, G., Witten, D., Hastie, T., Tibshirani, R. & Taylor, J. (2023), *An Introduction to Statistical Learning with Applications in Python*, Springer Texts in Statistics, Springer, Cham.

**URL:** <https://link.springer.com/book/10.1007/978-3-031-38747-0>

Kahl, S., Wood, C. M., Eibl, M. & Klinck, H. (2021), 'Birdnet: A deep learning solution for avian diversity monitoring', *Ecological Informatics* **61**, 101236.

**URL:** <https://www.sciencedirect.com/science/article/pii/S1574954121000273>

Keren, G. & Schuller, B. (2016), 'Convolutional rnn: An enhanced model for extracting features from sequential data', pp. 3412-3419.

- Kolecki, J. C. (2002), ‘An introduction to tensors for students of physics and engineering’.
- Krizhevsky, A., Sutskever, I. & Hinton, G. E. (2017), ‘Imagenet classification with deep convolutional neural networks’, *Commun. ACM* **60**(6), 84–90.  
**URL:** <https://doi.org/10.1145/3065386>
- Lecun, Y., Bottou, L., Bengio, Y. & Haffner, P. (1998), ‘Gradient-based learning applied to document recognition’, *Proceedings of the IEEE* **86**(11), 2278–2324.  
**URL:** <http://ieeexplore.ieee.org/document/726791/>
- Li, F.-F. & Adeli, E. (2024), ‘CS231n: Deep Learning for Computer Vision’.  
**URL:** <https://cs231n.stanford.edu/>
- Lu, T., Han, B. & Yu, F. (2021), ‘Detection and classification of marine mammal sounds using alexnet with transfer learning’, *Ecological Informatics* **62**, 101277.  
**URL:** <https://www.sciencedirect.com/science/article/pii/S1574954121000686>
- Ma, L., Milner, B. & Smith, D. (2006), ‘Acoustic environment classification’, *ACM Transactions on Speech and Language Processing* **3**(2), 1–22.  
**URL:** <https://dl.acm.org/doi/10.1145/1149290.1149292>
- McCloskey, E. (2024), ‘Interview with the school of biology on false killer whales audio’.
- Mutanu, L., Gohil, J., Gupta, K., Wagio, P. & Kotonya, G. (2022), ‘A Review of Automated Bioacoustics and General Acoustics Classification Research’, *Sensors* **22**(21), 8361.  
**URL:** <https://www.mdpi.com/1424-8220/22/21/8361>
- National Oceanography Centre (2024), ‘Gliders’.
- NOAA Fisheries (2024), ‘Mission on the high seas: Hawaiian islands cetacean and ecosystem assessment survey’.  
**URL:** <https://www.fisheries.noaa.gov/feature-story/mission-high-seas-hawaiian-islands-cetacean-and-ecosystem-assessment-survey>
- Nur Korkmaz, B., Diamant, R., Danino, G. & Testolin, A. (2023), ‘Automated

detection of dolphin whistles with convolutional networks and transfer learning’, *Frontiers in Artificial Intelligence* **6**.

**URL:** <https://www.frontiersin.org/journals/artificial-intelligence/articles/10.3389/frai.2023.1099022>

Oppenheim, A. V., Schafer, R. W. & Buck, J. R. (1998), *Discrete-time signal processing*, 2 edn, Pearson, Upper Saddle River, NJ.

O’Shea, K. & Nash, R. (2015), ‘An introduction to convolutional neural networks’.

**URL:** <https://arxiv.org/abs/1511.08458>

Palanisamy, K., Singhanian, D. & Yao, A. (2020), ‘Rethinking CNN models for audio classification’.

**URL:** <http://arxiv.org/abs/2007.11154>

*Papers With Code: Trends* (2024).

**URL:** <https://paperswithcode.com/trends>

Parsons, S. (2001), ‘Identification of New Zealand bats ( *Chalinolobus tuberculatus* and *Mystacina tuberculata* ) in flight from analysis of echolocation calls by artificial neural networks’, *Journal of Zoology* **253**(4), 447–456.

**URL:** <https://zslpublications.onlinelibrary.wiley.com/doi/10.1017/S0952836901000413>

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., KÄpf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J. & Chintala, S. (2019), ‘Pytorch: An imperative style, high-performance deep learning library’.

**URL:** <https://arxiv.org/abs/1912.01703>

Phung, V. H. & Rhee, E. J. (2019), ‘A High-Accuracy Model Average Ensemble of Convolutional Neural Networks for Classification of Cloud Image Patches on Small Datasets’, *Applied Sciences* **9**(21), 4500.

**URL:** <https://www.mdpi.com/2076-3417/9/21/4500>

Pixabay (2016), ‘Puppy, dog, pet image. free for use.’.

**URL:** <https://pixabay.com/photos/puppy-dog-pet-collar-dog-collar-1903313/>

*Pseudorca crassidens*: Baird, R.W.: *The IUCN Red List of Threatened Species 2018: e.T18596A145357488* (2018).

**URL:** <http://dx.doi.org/10.2305/IUCN.UK.2018-2.RLTS.T18596A145357488.en>

Pytorch (2024), ‘Tensors’.

**URL:** [https://pytorch.org/tutorials/beginner/basics/tensorqs\\_tutorial.html](https://pytorch.org/tutorials/beginner/basics/tensorqs_tutorial.html)

Razavian, A. S., Azizpour, H., Sullivan, J. & Carlsson, S. (2014), ‘CNN Features off-the-shelf: an Astounding Baseline for Recognition’. arXiv:1403.6382 [cs].

**URL:** <http://arxiv.org/abs/1403.6382>

Ridgway, S. H. & Harrison, R. J. (1999), *Handbook of marine mammals*, Academic press.

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C. & Fei-Fei, L. (2015), ‘Imagenet large scale visual recognition challenge’.

**URL:** <https://arxiv.org/abs/1409.0575>

Sayigh, L., Daher, M. A., Allen, J., Gordon, H., Joyce, K., Stuhlmann, C. & Tyack, P. (2017), ‘The Watkins Marine Mammal Sound Database: An online, freely accessible resource’, *Proceedings of Meetings on Acoustics* **27**(1), 040013.

**URL:** <https://doi.org/10.1121/2.0000358>

Simonyan, K. & Zisserman, A. (2015), ‘Very Deep Convolutional Networks for Large-Scale Image Recognition’. arXiv:1409.1556 [cs].

**URL:** <http://arxiv.org/abs/1409.1556>

Song, G., Wang, Z., Han, F. & Ding, S. (2017), Transfer Learning for Music Genre Classification, in Z. Shi, B. Goertzel & J. Feng, eds, ‘Intelligence Science I’, Vol. 510, Springer International Publishing, Cham, pp. 183–190. Series Title: IFIP Advances in Information and Communication Technology.

**URL:** [https://link.springer.com/10.1007/978-3-319-68121-4\\_19](https://link.springer.com/10.1007/978-3-319-68121-4_19)

Sousa-Lima, R. (2013), ‘Errata: Aquatic mammals, 39(1), 2013, pp. 23-53 a review and inventory of fixed autonomous recorders for passive acoustic monitoring of marine



mammals’, *Aquatic Mammals* **39**(2), 205-210.

**URL:** <http://dx.doi.org/10.1578/AM.39.2.2013.205>

Springenberg, J. T., Dosovitskiy, A., Brox, T. & Riedmiller, M. (2015), ‘Striving for simplicity: The all convolutional net’.

**URL:** <https://arxiv.org/abs/1412.6806>

Sturm, B. L. (2014), ‘The GTZAN dataset: Its contents, its faults, their effects on evaluation, and its future use’, *Journal of New Music Research* **43**(2), 147–172. arXiv:1306.1461 [cs].

**URL:** <http://arxiv.org/abs/1306.1461>

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V. & Rabinovich, A. (2014), ‘Going Deeper with Convolutions’. arXiv:1409.4842 [cs].

**URL:** <http://arxiv.org/abs/1409.4842>

Tan, M. & Le, Q. V. (2020), ‘EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks’. arXiv:1905.11946 [cs, stat].

**URL:** <http://arxiv.org/abs/1905.11946>

*torch.utils.data* (2023).

**URL:** <https://pytorch.org/docs/stable/data.html>

*Torchvision* (2024).

**URL:** <https://pytorch.org/vision/stable/index.html>

Williams, B., van Merriënboer, B., Dumoulin, V., Hamer, J., Triantafillou, E., Fleishman, A. B., McKown, M., Munger, J. E., Rice, A. N., Lillis, A., White, C. E., Hobbs, C. A. D., Razak, T. B., Jones, K. E. & Denton, T. (2024), ‘Leveraging tropical reef, bird and unrelated sounds for superior transfer learning in marine bioacoustics’.

**URL:** <https://arxiv.org/abs/2404.16436>

Yadav, S. & Bhole, G. P. (2020), Handling imbalanced dataset classification in machine

learning, in ‘2020 IEEE Pune Section International Conference (PuneCon)’, pp. 38–43.

Yani, M., Irawan, S. & Setianingsih, C. (2019), ‘Application of transfer learning using convolutional neural network method for early detection of terry’s nail’, *Journal of Physics: Conference Series* **1201**, 012052.

Yano, K. M., Oleson, E. M. E. M., Keating, J. L., Ballance, L. T. L. T., Hill, M. C. M. C., Bradford, A. L., Allen, A. N., Joyce, T. W., Moore, J. E. & Henry, A. E. (2018), ‘Cetacean and seabird data collected during the hawaiian islands cetacean and ecosystem assessment survey (HICEAS), july - december 2017’. Publisher: United States. National Marine Fisheries Service Pacific Islands Fisheries Science Center (U.S.) BOEM Alaska Environmental Studies Program.

**URL:** <https://repository.library.noaa.gov/view/noaa/18660>

Zhang, A., Lipton, Z. C., Li, M. & Smola, A. J. (2023), *Dive into Deep Learning*, Cambridge University Press. <https://D2L.ai>.

Zhong, M., LeBien, J., Campos-Cerqueira, M., Dodhia, R., Lavista Ferres, J., Velez, J. P. & Aide, T. M. (2020), ‘Multispecies bioacoustic classification using transfer learning of deep convolutional neural networks with pseudo-labeling’, *Applied Acoustics* **166**, 107375.

**URL:** <https://linkinghub.elsevier.com/retrieve/pii/S0003682X20304795>

## A | Python Code

This code is for my MSc Data-Intensive Analysis dissertation entitled "*Convolutional Neural Networks for the Classification of Pseudorca crassidens Passive Acoustic Data*".

The version for this dissertation can be found in

<https://github.com/josh-arrabaca/standrewsdissertation/releases/tag/Submission3>.

To run these, please ensure the code is in the top folder, and the related dataset is in the subfolder "data", and that the audio files are in the appropriate sub-subfolders "Insular" and "Pelagic".

To make predictions, you can run at the command line: "python c\_make\_prediction.py filename.wav".

The flow of experiments is as follows:

1. For converting the audio files to spectrograms:
  - a\_convert\_wav\_to\_png.py
2. For applying Transfer Learning to different pretrained CNNs. These mostly use the same code, except a different model was loaded for each experiment:
  - b\_CNN\_model\_DenseNet.py
  - b\_CNN\_model\_googlenet.py
  - b\_CNN\_model\_efficientnet.py
  - b\_CNN\_model\_ResNet.py
3. For applying the *Weighed Random Sampler* to the imbalanced dataset. Much of the code is similar to the above, except for adding the sampler library, and plotting the new distribution per batch:
  - b\_CNN\_model\_ResNet\_wrs.py

4. For hyperparameter tuning. Tuning was done first with the learning rate, then the momentum value. The code is also very similar to the above, except for the for-loop code added to train and test each of the hyperparameter values:

- `b_CNN_model_ResNet_wrs_lr.py`
- `b_CNN_model_ResNet_wrs_momentum.py`

5. For testing the best model with the unseen test set. As with the above, most of the code is the same, except for the metrics:

- `b_CNN_model_ResNet_best.py`

6. For making predictions using the saved best model:

- `c_make_prediction.py`

Code Listing A.1: code/a\_convert\_wav\_to\_png.py

```
1 # Import Libraries
2 import os.path
3 import numpy as np
4 from pydub import AudioSegment
5 from scipy import signal
6 import matplotlib.pyplot as plt
7
8 # Define the function that will convert the files
9 def batch_wav_to_png ( folder ):
10     """ Goes through all subfolders in a chosen folder , and converts wav files to
11         spectrogram png. """
12     print("Finding files and converting wav to png...\n")
13
14     # Create folders for the images
15     if not os.path.exists(os.path.join("data", "images", "Insular")):
16         os.makedirs(os.path.join("data", "images", "Insular"))
17     if not os.path.exists(os.path.join("data", "images", "Pelagic")):
18         os.makedirs(os.path.join("data", "images", "Pelagic"))
19
20     # Initialise counters
21     wavfile_counter = 0
22     otherfile_counter = 0
23     folder_counter = 0
24
25     # Walkthrough all subfolders
26     for root, dirs, files in os.walk(folder):
27
28         for file in files:
29
30             # Process files if they are wav
31             if file.endswith((".wav")) and any(dolphin_type in str(root) for
32 dolphin_type in ["Pelagic", "Insular"]):
33                 fileloc = os.path.join(root, file)
34
35                 # Extract raw audio data from the wav, and "flatten" channels to only
36                 one channel
37                 sound = AudioSegment.from_wav( fileloc )
38                 sound = sound.set_channels(1)
39
40                 # Convert the result to ndarray, and find the sample_rate
41                 samples = sound.get_array_of_samples()
42                 samples = np.array(samples).astype(np.int16)
43                 sample_rate = sound.frame_rate
```

```
42         # Convert to spectrogram data
43         frequencies , times , spectrogram = signal.spectrogram(samples ,
sample_rate)
44
45         # Plot the result without any axes or labels
46         # Part of this code was provided by Dr Chrissy Fell , School of
Mathematics and Statistics , University of St Andrews
47         plt.pcolormesh(times , frequencies , np.log(spectrogram))
48         plt.tick_params(axis='both' , which='both' , bottom=False , left=False ,
top=False ,
49                             labelbottom=False , labelleft=False)
50         plt.rcParams["figure.figsize"] = (12,8)
51         if "Insular" in str(root):
52             my_file = file + ".png"
53             plt.savefig(os.path.join(folder , "images" , "Insular" , my_file))
54         elif "Pelagic" in str(root):
55             my_file = file + ".png"
56             plt.savefig(os.path.join(folder , "images" , "Pelagic" , my_file))
57         plt.close()
58
59         #Increment the counters
60         wavfile_counter += 1
61         if wavfile_counter % 50 == 0:
62             print (f"Processing {wavfile_counter}th wave file ...")
63
64         else:
65             otherfile_counter += 1
66
67         folder_counter += 1
68
69         print (f"Processing done!\nConverted {wavfile_counter} .wav files to .png in {
folder_counter} folders.\nThe files are located in '{folder}\images' folder.\n\
nIgnored {otherfile_counter} other files.")
70
71
72 # Name the folder with wav files here
73 folder = "data"
74
75 # Run the function
76 batch_wav_to_png (folder)
```

## Code Listing A.2: code/b\_CNN\_model\_DenseNet.py

```
1 # Adapted from the original Pytorch tutorial by Sasank Chilamkurthy
2
3 # Import libraries
4 import torch
5 import torch.nn as nn
6 import torch.optim as optim
7 from torch.optim import lr_scheduler
8 import torch.backends.cudnn as cudnn
9 import torchvision
10 from torchvision import datasets, models, transforms
11 import matplotlib.pyplot as plt
12 import time
13 import os
14 from PIL import Image
15 from tempfile import TemporaryDirectory
16 from torch.utils.data import DataLoader
17 from sklearn import metrics
18
19
20 cudnn.benchmark = True
21 plt.ion() # interactive mode
22
23 # Set 'random' seed
24 torch.manual_seed(220029955)
25
26 # Welcome message
27 print("Welcome! We will train the last layer of a pre-trained CNN model.\n")
28
29
30 # Define the transforms needed
31 data_transforms = transforms.Compose([
32     transforms.Resize([224,224]), # Minimum size needed for Densenet
33     transforms.ToTensor(),
34     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
35     normalisation for Densenet
36 ])
37
38 # Get the dataset from the images created from the wav files
39 dataset = datasets.ImageFolder(os.path.join("data", "images"), transform=
40     data_transforms)
41
42 # Define the classes (Insular and Pelagic)
43 classes = dataset.classes
```

```
43 # Split the data into train, val and test sets
44 train_size = int(0.6 * len(dataset))
45 val_size = int((len(dataset) - train_size) / 2)
46 test_size = val_size
47 train_dataset, val_dataset, test_dataset = torch.utils.data.random_split(dataset, [
    train_size, val_size, test_size])
48 print(f"The dataset consists of {train_size + val_size + test_size} datapoints, split
    as follows:")
49 print(f"Train set: {train_size} \nValidation set: {val_size} \nTest size: {test_size}
    \n")
50
51 # Define the device to be used for training
52 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
53 print(f'The device being used is: {device}\n')
54
55 # Define the batch size and number of epochs based on the device
56 if str(device) == "cuda:0":
57     batch_size = 64
58     num_epochs = 24
59 else:
60     batch_size = 20
61     num_epochs = 3
62
63 # Dataloaders
64 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
65 val_dataloader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
66 test_dataloader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
67
68 dataloaders = {"train": train_dataloader,
69               "val": val_dataloader}
70
71 dataset_sizes = {"train": len(train_dataset),
72                 "val": len(test_dataset)}
73
74
75 # below code is taken from
76 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
77
78     since = time.time()
79
80     # Create a temporary directory to save training checkpoints
81     with TemporaryDirectory() as tempdir:
82         best_model_params_path = os.path.join(tempdir, 'best_model_params.pt')
83
84         torch.save(model.state_dict(), best_model_params_path)
85         best_acc = 0.0
```



```
86
87     for epoch in range(num_epochs):
88         print(f'Epoch {epoch}/{num_epochs - 1}')
89         print('-' * 10)
90
91         # Each epoch has a training and validation phase
92         for phase in ['train', 'val']:
93             if phase == 'train':
94                 model.train() # Set model to training mode
95             else:
96                 model.eval() # Set model to evaluate mode
97
98             running_loss = 0.0
99             running_corrects = 0
100
101             # Iterate over data.
102             for inputs, labels in dataloaders[phase]:
103                 inputs = inputs.to(device)
104                 labels = labels.to(device)
105
106                 # zero the parameter gradients
107                 optimizer.zero_grad()
108
109                 # forward
110                 # track history if only in train
111                 with torch.set_grad_enabled(phase == 'train'):
112                     outputs = model(inputs)
113                     _, preds = torch.max(outputs, 1)
114                     loss = criterion(outputs, labels)
115
116                 # backward + optimize only if in training phase
117                 if phase == 'train':
118                     loss.backward()
119                     optimizer.step()
120
121                 # statistics
122                 running_loss += loss.item() * inputs.size(0)
123                 running_corrects += torch.sum(preds == labels.data)
124             if phase == 'train':
125                 scheduler.step()
126
127             epoch_loss = running_loss / dataset_sizes[phase]
128             epoch_acc = running_corrects.double() / dataset_sizes[phase]
129
130             print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
131
```

```
132         # deep copy the model
133         if phase == 'val' and epoch_acc > best_acc:
134             best_acc = epoch_acc
135             torch.save(model.state_dict(), best_model_params_path)
136
137         print()
138
139         time_elapsed = time.time() - since
140         print(f'Training complete in {time_elapsed // 60:.0f}m {time_elapsed % 60:.0f}
141         s\n\n')
142         print(f'The best val accuracy score is: {best_acc:4f}\n\n')
143
144         # load best model weights
145         model.load_state_dict(torch.load(best_model_params_path))
146     return model
147
148 # Load Densenet
149 model_conv = torchvision.models.densenet121(weights='IMAGENET1K_V1')
150
151 # This part does the training on the final layer only
152 for param in model_conv.parameters():
153     param.requires_grad = False
154
155 # Parameters of newly constructed modules have requires_grad=True by default
156 num_ftrs = model_conv.classifier.in_features
157 model_conv.classifier = nn.Linear(num_ftrs, 2)
158
159 model_conv = model_conv.to(device)
160
161 criterion = nn.CrossEntropyLoss()
162
163 # Observe that only parameters of final layer are being optimized as
164 # opposed to before.
165 optimizer_conv = optim.SGD(model_conv.classifier.parameters(), lr=0.001, momentum
166                             =0.9)
167
168 # Decay LR by a factor of 0.1 every 7 epochs
169 exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
170
171 # Now train the model, and view the loss and accuracy scores
172 model_conv = train_model(model_conv, criterion, optimizer_conv,
173                           exp_lr_scheduler, num_epochs=num_epochs)
174
175 # This part creates the classes and scores from the validation data for the metrics
```

```
176 y_score = []
177 true_classes = []
178 predicted_classes = []
179
180 for inputs, labels in val_dataloader:
181     inputs = inputs.to(device)
182     labels = labels.to(device)
183
184     # This part 'flattens' the tensor into a list
185     labels = labels.cpu().numpy().tolist()
186     true_classes.extend(labels)
187
188     with torch.no_grad():
189         model_conv.eval()
190         output = model_conv(inputs)
191
192         for each_output in output:
193             predicted_class = each_output.cpu().data.numpy().argmax() # Numpify each
194             predicted_classes.append(predicted_class) # Conactenate
195
196             p = torch.nn.functional.softmax(output, dim=1) # Get probabilities
197             top_proba = p.cpu().numpy()[0][0] # Get probabilities of the positive
198             class ('Insular') only
199             y_score.append(top_proba)# predicted_classes
200
201 # Let's find the precision, recall and confusion matrix as well
202 print(f"Precision: {round(metrics.precision_score(true_classes, predicted_classes),5)}")
203 print(f"Recall: {round(metrics.recall_score(true_classes, predicted_classes),5)}")
204 print("Confusion matrix:\n", metrics.confusion_matrix(true_classes, predicted_classes))
```

## Code Listing A.3: code/b\_CNN\_model\_efficientnet.py

```
1 # Adapted from the original Pytorch tutorial by Sasank Chilamkurthy
2
3 # Import libraries
4 import torch
5 import torch.nn as nn
6 import torch.optim as optim
7 from torch.optim import lr_scheduler
8 import torch.backends.cudnn as cudnn
9 import torchvision
10 from torchvision import datasets, models, transforms
11 import matplotlib.pyplot as plt
12 import time
13 import os
14 from PIL import Image
15 from tempfile import TemporaryDirectory
16 from torch.utils.data import DataLoader
17 from sklearn import metrics
18
19 cudnn.benchmark = True
20 plt.ion() # interactive mode
21
22 # Set 'random' seed
23 torch.manual_seed(220029955)
24
25 # Welcome message
26 print("Welcome! We will train the last layer of a pre-trained CNN model.\n")
27
28 # Define the transforms needed
29 data_transforms = transforms.Compose([
30     transforms.Resize([384,384]), # Minimum size needed for Efficientnet
31     transforms.ToTensor(),
32     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
33     normalisation for Efficientnet
34 ])
35
36 # Get the dataset from the images created from the wav files
37 dataset = datasets.ImageFolder(os.path.join("data", "images"), transform=
38     data_transforms)
39
40 # Define the classes (Insular and Pelagic)
41 classes = dataset.classes
42
43 # Split the data into train, val and test sets
44 train_size = int(0.6 * len(dataset))
```

```
43 val_size = int((len(dataset) - train_size) / 2)
44 test_size = val_size
45 train_dataset, val_dataset, test_dataset = torch.utils.data.random_split(dataset, [
    train_size, val_size, test_size])
46 print(f"The dataset consists of {train_size + val_size + test_size} datapoints, split
    as follows:")
47 print(f"Train set: {train_size} \nValidation set: {val_size} \nTest size: {test_size}
    \n")
48
49 # Define the device to be used for training
50 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
51 print(f'The device being used is: {device}\n')
52
53 # Define the batch size and number of epochs based on the device
54 if str(device) == "cuda:0":
55     batch_size = 64
56     num_epochs = 24
57 else:
58     batch_size = 20
59     num_epochs = 3
60
61 # Dataloaders
62 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
63 val_dataloader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
64 test_dataloader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
65
66 dataloaders = {"train": train_dataloader,
67               "val": val_dataloader}
68
69 dataset_sizes = {"train": len(train_dataset),
70                 "val": len(test_dataset)}
71
72
73 # below code is taken from
74 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
75
76     since = time.time()
77
78     # Create a temporary directory to save training checkpoints
79     with TemporaryDirectory() as tempdir:
80         best_model_params_path = os.path.join(tempdir, 'best_model_params.pt')
81
82         torch.save(model.state_dict(), best_model_params_path)
83         best_acc = 0.0
84
85         for epoch in range(num_epochs):
```

```
86     print(f'Epoch {epoch}/{num_epochs - 1}')
87     print('-' * 10)
88
89     # Each epoch has a training and validation phase
90     for phase in ['train', 'val']:
91         if phase == 'train':
92             model.train() # Set model to training mode
93         else:
94             model.eval() # Set model to evaluate mode
95
96         running_loss = 0.0
97         running_corrects = 0
98
99         # Iterate over data.
100        for inputs, labels in dataloaders[phase]:
101            inputs = inputs.to(device)
102            labels = labels.to(device)
103
104            # zero the parameter gradients
105            optimizer.zero_grad()
106
107            # forward
108            # track history if only in train
109            with torch.set_grad_enabled(phase == 'train'):
110                outputs = model(inputs)
111                _, preds = torch.max(outputs, 1)
112                loss = criterion(outputs, labels)
113
114            # backward + optimize only if in training phase
115            if phase == 'train':
116                loss.backward()
117                optimizer.step()
118
119            # statistics
120            running_loss += loss.item() * inputs.size(0)
121            running_corrects += torch.sum(preds == labels.data)
122        if phase == 'train':
123            scheduler.step()
124
125        epoch_loss = running_loss / dataset_sizes[phase]
126        epoch_acc = running_corrects.double() / dataset_sizes[phase]
127
128        print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
129
130        # deep copy the model
131        if phase == 'val' and epoch_acc > best_acc:
```

```
132         best_acc = epoch_acc
133         torch.save(model.state_dict(), best_model_params_path)
134
135         print()
136
137         time_elapsed = time.time() - since
138         print(f'Training complete in {time_elapsed // 60:.0f}m {time_elapsed % 60:.0f}
139         }s\n\n')
140         print(f'The best val accuracy score is: {best_acc:4f}\n\n')
141
142         # load best model weights
143         model.load_state_dict(torch.load(best_model_params_path))
144     return model
145
146 # Loading efficientnet
147 model_conv = torchvision.models.efficientnet_v2_s(weights='IMAGENET1K_V1')
148
149
150 # This part does the training on the final layer only
151 for param in model_conv.parameters():
152     param.requires_grad = False
153
154 # Parameters of newly constructed modules have requires_grad=True by default
155 num_fts = model_conv.classifier[1].in_features
156 model_conv.classifier[1] = nn.Linear(num_fts, 2)
157
158 model_conv = model_conv.to(device)
159
160 criterion = nn.CrossEntropyLoss()
161
162 # Observe that only parameters of final layer are being optimized as
163 # opposed to before.
164 optimizer_conv = optim.SGD(model_conv.classifier[1].parameters(), lr=0.001, momentum
165                             =0.9)
166
167 # Decay LR by a factor of 0.1 every 7 epochs
168 exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
169
170 # Now train the model, and view the loss and accuracy scores
171 model_conv = train_model(model_conv, criterion, optimizer_conv,
172                           exp_lr_scheduler, num_epochs=num_epochs)
173
174 # This part creates the classes and scores from the validation data for the metrics
175 y_score = []
176 true_classes = []
```

```
176 predicted_classes = []
177
178 for inputs, labels in val_dataloader:
179     inputs = inputs.to(device)
180     labels = labels.to(device)
181
182     # This part 'flattens' the tensor into a list
183     labels = labels.cpu().numpy().tolist()
184     true_classes.extend(labels)
185
186     with torch.no_grad():
187         model_conv.eval()
188         output = model_conv(inputs)
189
190         for each_output in output:
191             predicted_class = each_output.cpu().data.numpy().argmax() # Numpify each
192             # output
193             predicted_classes.append(predicted_class) # Conactenate
194
195             p = torch.nn.functional.softmax(output, dim=1) # Get probabilities
196             top_proba = p.cpu().numpy()[0][0] # Get probabilities of the positive
197             # class ('Insular') only
198             y_score.append(top_proba)# predicted_classes
199
200 # Let's find the precision, recall and confusion matrix as well
201 print(f"Precision: {round(metrics.precision_score(true_classes, predicted_classes),5)}")
202 print(f"Recall: {round(metrics.recall_score(true_classes, predicted_classes),5)}")
203 print("Confusion matrix:\n", metrics.confusion_matrix(true_classes, predicted_classes))
```



Code Listing A.4: code/b\_CNN\_model\_ResNet.py

```
1 # Adapted from the original Pytorch tutorial by Sasank Chilamkurthy
2
3 # Import libraries
4 import torch
5 import torch.nn as nn
6 import torch.optim as optim
7 from torch.optim import lr_scheduler
8 import torch.backends.cudnn as cudnn
9 import numpy as np
10 import torchvision
11 from torchvision import datasets, models, transforms
12 import matplotlib.pyplot as plt
13 import time
14 import os
15 from PIL import Image
16 from tempfile import TemporaryDirectory
17 from torch.utils.data import DataLoader
18 from sklearn import metrics
19 from sklearn.metrics import PrecisionRecallDisplay
20 from sklearn.metrics import precision_recall_curve
21
22 cudnn.benchmark = True
23 plt.ion() # interactive mode
24
25 # Set 'random' seed
26 torch.manual_seed(220029955)
27
28 # Welcome message
29 print("Welcome! We will train the last layer of a pre-trained CNN model.\n")
30
31
32 # Define the transforms needed
33 data_transforms = transforms.Compose([
34     transforms.Resize([224,224]), # Minimum size needed for Resnet
35     transforms.ToTensor(),
36     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
37     normalisation for Resnet
38 ])
39
40 # Get the dataset from the images created from the wav files
41 dataset = datasets.ImageFolder(os.path.join("data", "images"), transform=
42     data_transforms)
43
44 # Define the classes (Insular and Pelagic)
```

```
43 classes = dataset.classes
44
45 # Split the data into train, val and test sets
46 train_size = int(0.6 * len(dataset))
47 val_size = int((len(dataset) - train_size) / 2)
48 test_size = val_size
49 train_dataset, val_dataset, test_dataset = torch.utils.data.random_split(dataset, [
    train_size, val_size, test_size])
50 print(f"The dataset consists of {train_size + val_size + test_size} datapoints, split
    as follows:")
51 print(f"Train set: {train_size} \nValidation set: {val_size} \nTest size: {test_size}
    \n")
52
53 # Define the device to be used for training
54 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
55 print(f'The device being used is: {device}\n')
56
57 # Define the batch size and number of epochs based on the device
58 if str(device) == "cuda:0":
59     batch_size = 64
60     num_epochs = 24
61 else:
62     batch_size = 20
63     num_epochs = 3
64
65 # Dataloaders
66 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
67 val_dataloader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
68 test_dataloader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
69
70 dataloaders = {"train": train_dataloader,
71               "val": val_dataloader}
72
73 dataset_sizes = {"train": len(train_dataset),
74                 "val": len(test_dataset)}
75
76
77 # below code is taken from
78 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
79
80     since = time.time()
81
82     # Create a temporary directory to save training checkpoints
83     with TemporaryDirectory() as tempdir:
84         best_model_params_path = os.path.join(tempdir, 'best_model_params.pt')
85
```

```
86     torch.save(model.state_dict(), best_model_params_path)
87     best_acc = 0.0
88
89     for epoch in range(num_epochs):
90         print(f'Epoch {epoch}/{num_epochs - 1}')
91         print('-' * 10)
92
93         # Each epoch has a training and validation phase
94         for phase in ['train', 'val']:
95             if phase == 'train':
96                 model.train() # Set model to training mode
97             else:
98                 model.eval() # Set model to evaluate mode
99
100             running_loss = 0.0
101             running_corrects = 0
102
103             # Iterate over data.
104             for inputs, labels in dataloaders[phase]:
105                 inputs = inputs.to(device)
106                 labels = labels.to(device)
107
108                 # zero the parameter gradients
109                 optimizer.zero_grad()
110
111                 # forward
112                 # track history if only in train
113                 with torch.set_grad_enabled(phase == 'train'):
114                     outputs = model(inputs)
115                     _, preds = torch.max(outputs, 1)
116                     loss = criterion(outputs, labels)
117
118                 # backward + optimize only if in training phase
119                 if phase == 'train':
120                     loss.backward()
121                     optimizer.step()
122
123                 # statistics
124                 running_loss += loss.item() * inputs.size(0)
125                 running_corrects += torch.sum(preds == labels.data)
126             if phase == 'train':
127                 scheduler.step()
128
129             epoch_loss = running_loss / dataset_sizes[phase]
130             epoch_acc = running_corrects.double() / dataset_sizes[phase]
131
```

```
132         print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
133
134         # deep copy the model
135         if phase == 'val' and epoch_acc > best_acc:
136             best_acc = epoch_acc
137             torch.save(model.state_dict(), best_model_params_path)
138
139         print()
140
141         time_elapsed = time.time() - since
142         print(f'Training complete in {time_elapsed // 60:.0f}m {time_elapsed % 60:.0f}
143         }s\n\n')
144         print(f'The best val accuracy score is: {best_acc:.4f}\n\n')
145
146         # load best model weights
147         model.load_state_dict(torch.load(best_model_params_path))
148     return model
149
150 # Loading ResNet
151 model_conv = torchvision.models.resnet18(weights='IMAGENET1K_V1')
152
153 # This part does the training on the final layer only
154 for param in model_conv.parameters():
155     param.requires_grad = False
156
157 # Parameters of newly constructed modules have requires_grad=True by default
158 num_fts = model_conv.fc.in_features
159 model_conv.fc = nn.Linear(num_fts, 2)
160
161 model_conv = model_conv.to(device)
162
163 criterion = nn.CrossEntropyLoss()
164
165 # Observe that only parameters of final layer are being optimized as
166 # opposed to before.
167 optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=0.001, momentum=0.9)
168
169 # Decay LR by a factor of 0.1 every 7 epochs
170 exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
171
172
173 # Now train the model, and view the loss and accuracy scores
174 model_conv = train_model(model_conv, criterion, optimizer_conv,
175                          exp_lr_scheduler, num_epochs=num_epochs)
176
```

```
177 # This part creates the classes and scores from the validation data for the plots
178 y_score = []
179 true_classes = []
180 predicted_classes = []
181
182 for inputs, labels in val_dataloader:
183     inputs = inputs.to(device)
184     labels = labels.to(device)
185
186     # This part 'flattens' the tensor into a list
187     labels = labels.cpu().numpy().tolist()
188     true_classes.extend(labels)
189
190     with torch.no_grad():
191         model_conv.eval()
192         output = model_conv(inputs)
193
194         for each_output in output:
195             predicted_class = each_output.cpu().data.numpy().argmax() # Numpify each
196             # output
197             predicted_classes.append(predicted_class) # Conactenate
198
199             p = torch.nn.functional.softmax(output, dim=1) # Get probabilities
200             top_proba = p.cpu().numpy()[0][0] # Get probabilities of the positive
201             # class ('Insular') only
202             y_score.append(top_proba) # predicted_classes
203
204 # Let's also check the Precision-Recall Curve, and the Threshold Plots
205 # These plots are saved locally as image files
206 display = PrecisionRecallDisplay.from_predictions(true_classes,
207                                                  y_score,
208                                                  name="CNN",
209                                                  plot_chance_level=True)
210 _ = display.ax_.set_title("2-class Precision-Recall curve")
211 plt.ylim([0,1])
212 plt.savefig("precision-recall-curve.png", bbox_inches='tight')
213 plt.close()
214
215 precision, recall, thresholds = precision_recall_curve(true_classes, y_score)
216 plt.plot(thresholds, precision[:-1], 'b--', label='Precision')
217 plt.plot(thresholds, recall[:-1], 'r--', label='Recall')
218 plt.xlabel('Threshold')
219 plt.legend(loc='lower left')
220 plt.ylim([0,1])
221 plt.title("Precision-Recall-Threshold Plot")
222 plt.savefig("precision-recall-threshold.png", bbox_inches='tight')
```

```
221
222 # Let's find the precision, recall and confusion matrix as well
223 print(f"Precision: {round(metrics.precision_score(true_classes, predicted_classes),5)}")
224 print(f"Recall: {round(metrics.recall_score(true_classes, predicted_classes),5)}")
225 print("Confusion matrix:\n", metrics.confusion_matrix(true_classes, predicted_classes))
```

Code Listing A.5: code/b\_CNN\_model\_ResNet\_wrs.py

```
1 # Import libraries
2 import torch
3 import torch.nn as nn
4 import torch.optim as optim
5 from torch.optim import lr_scheduler
6 import torch.backends.cudnn as cudnn
7 import numpy as np
8 import torchvision
9 from torchvision import datasets, models, transforms
10 import matplotlib.pyplot as plt
11 import time
12 import os
13 from PIL import Image
14 from tempfile import TemporaryDirectory
15 from torch.utils.data import DataLoader
16 from sklearn.metrics import PrecisionRecallDisplay
17 from sklearn.metrics import precision_recall_curve
18 from torch.utils.data import DataLoader
19 import pandas as pd
20 from torch.utils.data import WeightedRandomSampler
21
22 cudnn.benchmark = True
23 plt.ion() # interactive mode
24
25 # Set 'random' seed
26 torch.manual_seed(220029955)
27
28 # Welcome message
29 print("Welcome! We will train the last layer of a pre-trained CNN model.\n")
30
31
32 # Define the transforms needed
33 data_transforms = transforms.Compose([
34     transforms.Resize([224,224]), # Minimum size needed for Densenet
35     transforms.ToTensor(),
36     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
37     normalisation for Densenet
38 ])
39
40 # Get the dataset from the images created from the wav files
41 dataset = datasets.ImageFolder(os.path.join("data", "images"), transform=
42     data_transforms)
43
44 # Define the classes (Insular and Pelagic)
```

```
43 classes = dataset.classes
44
45 # Split the data into train, val and test sets
46 train_size = int(0.6 * len(dataset))
47 val_size = int((len(dataset) - train_size) / 2)
48 test_size = val_size
49 train_dataset, val_dataset, test_dataset = torch.utils.data.random_split(dataset, [
    train_size, val_size, test_size])
50 print(f"The dataset consists of {train_size + val_size + test_size} datapoints, split
    as follows:")
51 print(f"Train set: {train_size} \nValidation set: {val_size} \nTest size: {test_size}
    \n")
52
53 # Define the device to be used for training
54 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
55 print(f'The device being used is: {device}\n')
56
57 # Define the batch size and number of epochs based on the device
58 if str(device) == "cuda:0":
59     batch_size = 64
60     num_epochs = 24
61 else:
62     batch_size = 20
63     num_epochs = 3
64
65 # Dataloaders
66 # train_dataloader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
67 val_dataloader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
68 test_dataloader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
69
70
71 # Define a function to visualise the class distributions per batch
72 def plot_classes(dataloader):
73     batch_num_counts = []
74     cur_batch = 0
75     class_0_counts = []
76     class_1_counts = []
77     for i in range(len(train_dataloader)):
78         my_inputs, my_labels = next(iter(train_dataloader))
79         this_class_0_count = 0
80         this_class_1_count = 0
81         for label in my_labels:
82             if label.cpu().numpy() == 0:
83                 this_class_0_count+=1
84             elif label.cpu().numpy() == 1:
85                 this_class_1_count+=1
```



```
86         class_0_counts.append(this_class_0_count)
87         class_1_counts.append(this_class_1_count)
88         batch_num_counts.append(cur_batch)
89         cur_batch += 1
90
91     df = pd.DataFrame({'Insular': class_0_counts,
92                       'Pelagic': class_1_counts})
93
94     df.plot(kind="bar",
95            title="Weighed Distribution of Classes per Batch",
96            xlabel="Batch Number",
97            ylabel="Images per Batch")
98
99     plt.savefig("weighed_distribution_pre_optim.png", bbox_inches='tight')
100    plt.close()
101    print(df)
102
103    ## Visualise the distribution
104    # plot_classes(train_dataloader, "Imbalanced_Distribution_of_Classes_per_Batch")
105
106    # Define function to create a balanced sampler
107    # The following code was adapted from Vivek Maskara
108    # https://www.maskaravivek.com/post/pytorch-weighted-random-sampler/
109    def balanced_sampler(full_dataset, train_dataset):
110        # Find number of samples per class
111        y_train_indices = train_dataset.indices
112        y_train = [full_dataset.targets[i] for i in y_train_indices]
113        class_sample_count = np.array([len(np.where(y_train == t)[0]) for t in np.unique(y_train)])
114
115        # Find weights per class
116        weight = 1. / class_sample_count
117        samples_weight = np.array([weight[t] for t in y_train])
118        samples_weight = torch.from_numpy(samples_weight)
119
120        # Define sampler
121        sampler = WeightedRandomSampler(samples_weight.type('torch.DoubleTensor'), len(samples_weight))
122
123        return sampler
124
125    # Create a balanced sampler
126    sampler = balanced_sampler(dataset, train_dataset)
127    train_dataloader = DataLoader(train_dataset, batch_size=batch_size, sampler=sampler)
128
129    # Visualise the balanced dataset per batch
```

```
130 plot_classes(train_dataloader)
131
132 dataloaders = {"train": train_dataloader,
133               "val": val_dataloader}
134
135 dataset_sizes = {"train": len(train_dataset),
136                 "val": len(test_dataset)}
137
138
139 # Adapted from the original Pytorch tutorial by Sasank Chilamkurthy
140 train_losses = []
141 val_losses = []
142 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
143
144     since = time.time()
145
146     # Create a temporary directory to save training checkpoints
147     with TemporaryDirectory() as tempdir:
148         best_model_params_path = os.path.join(tempdir, 'best_model_params.pt')
149
150         torch.save(model.state_dict(), best_model_params_path)
151         best_acc = 0.0
152
153         for epoch in range(num_epochs):
154             print(f'Epoch {epoch}/{num_epochs - 1}')
155             print('-' * 10)
156
157             # Each epoch has a training and validation phase
158             for phase in ['train', 'val']:
159                 if phase == 'train':
160                     model.train() # Set model to training mode
161                 else:
162                     model.eval() # Set model to evaluate mode
163
164             running_loss = 0.0
165             running_corrects = 0
166
167             # Iterate over data.
168             for inputs, labels in dataloaders[phase]:
169                 inputs = inputs.to(device)
170                 labels = labels.to(device)
171
172                 # zero the parameter gradients
173                 optimizer.zero_grad()
174
175                 # forward
```

```
176         # track history if only in train
177         with torch.set_grad_enabled(phase == 'train'):
178             outputs = model(inputs)
179             _, preds = torch.max(outputs, 1)
180             loss = criterion(outputs, labels)
181
182         # backward + optimize only if in training phase
183         if phase == 'train':
184             loss.backward()
185             optimizer.step()
186
187         # statistics
188         running_loss += loss.item() * inputs.size(0)
189         running_corrects += torch.sum(preds == labels.data)
190     if phase == 'train':
191         scheduler.step()
192
193     epoch_loss = running_loss / dataset_sizes[phase]
194     epoch_acc = running_corrects.double() / dataset_sizes[phase]
195
196     if phase == 'train':
197         train_losses.append(epoch_loss)
198     else:
199         val_losses.append(epoch_loss)
200     print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
201
202
203
204     # deep copy the model
205     if phase == 'val' and epoch_acc > best_acc:
206         best_acc = epoch_acc
207         torch.save(model.state_dict(), best_model_params_path)
208
209     print()
210
211     time_elapsed = time.time() - since
212     print(f'Training complete in {time_elapsed // 60:.0f}m {time_elapsed % 60:.0f}
213     s\n\n')
214     print(f'The best val accuracy score is: {best_acc:.4f}\n\n')
215
216     # load best model weights
217     model.load_state_dict(torch.load(best_model_params_path))
218     return model
219
220 # We load ResNet since this had the best base accuracy compared to the other CNNs
```

```
221 model_conv = torchvision.models.resnet18(weights='IMAGENET1K_V1')
222
223 # This part does the training on the final layer only
224 for param in model_conv.parameters():
225     param.requires_grad = False
226
227 # Parameters of newly constructed modules have requires_grad=True by default
228 num_fts = model_conv.fc.in_features
229 model_conv.fc = nn.Linear(num_fts, 2)
230
231 model_conv = model_conv.to(device)
232
233 criterion = nn.CrossEntropyLoss()
234
235 # Observe that only parameters of final layer are being optimized as
236 # opposed to before.
237 optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=0.001, momentum=0.9)
238
239 # Decay LR by a factor of 0.1 every 7 epochs
240 exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
241
242 # Now we train the model
243 model_conv = train_model(model_conv, criterion, optimizer_conv,
244                           exp_lr_scheduler, num_epochs=num_epochs)
245
246
247 # This part creates the classes and scores from the training data for the plots
248 plt.plot(train_losses, 'b', label='Training Loss')
249 plt.plot(val_losses, 'r--', label='Validation Loss')
250 plt.legend(loc='upper right')
251 plt.xlabel('Epoch')
252 plt.ylabel('Loss')
253 plt.title("Model Loss over Epochs")
254 plt.savefig("losses.png", bbox_inches='tight')
255 plt.close()
256
257
258 y_score = []
259 true_classes = []
260 predicted_classes = []
261
262 for inputs, labels in train_dataloader:
263     inputs = inputs.to(device)
264     labels = labels.to(device)
265
266     # This part 'flattens' the tensor into a list
```

```
267 labels = labels.cpu().numpy().tolist()
268 true_classes.extend(labels)
269
270 with torch.no_grad():
271     model_conv.eval()
272     output = model_conv(inputs)
273
274     for each_output in output:
275         predicted_class = each_output.cpu().data.numpy().argmax() # Numpify each
276         predicted_classes.append(predicted_class) # Conactenate
277
278         p = torch.nn.functional.softmax(output, dim=1) # Get probabilities
279         top_proba = p.cpu().numpy()[0][0] # Get probabilities of the positive
280         y_score.append(top_proba)# predicted_classes
281
282 # Let's also check the Precision-Recall Curve, and the Threshold Plots after using
283 # the weighed sampler
284 # These plots are saved locally as image files
285 display = PrecisionRecallDisplay.from_predictions(true_classes,
286                                                  y_score,
287                                                  name="CNN",
288                                                  plot_chance_level=True)
289 _ = display.ax_.set_title("2-class Precision-Recall curve (Weighed Distribution)")
290 plt.ylim([0,1])
291 plt.savefig("precision-recall-curve-weighed.png", bbox_inches='tight')
292 plt.close()
293
294 precision, recall, thresholds = precision_recall_curve(true_classes, y_score)
295 plt.plot(thresholds, precision[:-1], 'b--', label='Precision')
296 plt.plot(thresholds, recall[:-1], 'r--', label='Recall')
297 plt.xlabel('Threshold')
298 plt.legend(loc='lower left')
299 plt.ylim([0,1])
300 plt.title("Precision-Recall-Threshold Plot (Weighed Distribution)")
301 plt.savefig("precision-recall-threshold-weighed.png", bbox_inches='tight')
302 plt.close()
```

Code Listing A.6: code/b\_CNN\_model\_ResNet\_wrs\_lr.py

```
1 # Import libraries
2 import torch
3 import torch.nn as nn
4 import torch.optim as optim
5 from torch.optim import lr_scheduler
6 import torch.backends.cudnn as cudnn
7 import numpy as np
8 import torchvision
9 from torchvision import datasets, models, transforms
10 import matplotlib.pyplot as plt
11 import time
12 import os
13 from PIL import Image
14 from tempfile import TemporaryDirectory
15 from torch.utils.data import DataLoader
16 import pandas as pd
17 from torch.utils.data import WeightedRandomSampler
18
19 cudnn.benchmark = True
20 plt.ion() # interactive mode
21
22 # Set 'random' seed
23 torch.manual_seed(220029955)
24
25 # Welcome message
26 print("Welcome! We will train the last layer of a pre-trained CNN model.\n")
27
28
29 # Define the transforms needed
30 data_transforms = transforms.Compose([
31     transforms.Resize([224,224]), # Minimum size needed for Densenet
32     transforms.ToTensor(),
33     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
34     normalisation for Densenet
35 ])
36
37 # Get the dataset from the images created from the wav files
38 dataset = datasets.ImageFolder(os.path.join("data", "images"), transform=
39     data_transforms)
40
41 # Define the classes (Insular and Pelagic)
42 classes = dataset.classes
```

```
43 train_size = int(0.6 * len(dataset))
44 val_size = int((len(dataset) - train_size) / 2)
45 test_size = val_size
46 train_dataset, val_dataset, test_dataset = torch.utils.data.random_split(dataset, [
    train_size, val_size, test_size])
47 print(f"The dataset consists of {train_size + val_size + test_size} datapoints, split
    as follows:")
48 print(f"Train set: {train_size} \nValidation set: {val_size} \nTest size: {test_size}
    \n")
49
50 # Define the device to be used for training
51 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
52 print(f'The device being used is: {device}\n')
53
54 # Define the batch size and number of epochs based on the device
55 if str(device) == "cuda:0":
56     batch_size = 64
57     num_epochs = 24
58 else:
59     batch_size = 20
60     num_epochs = 3
61
62 # Dataloaders
63 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
64 val_dataloader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
65 test_dataloader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
66
67
68 # Define function to create a balanced sampler
69 # https://www.maskaravivek.com/post/pytorch-weighted-random-sampler/
70 def balanced_sampler(full_dataset, train_dataset):
71     # Find number of samples per class
72     y_train_indices = train_dataset.indices
73     y_train = [full_dataset.targets[i] for i in y_train_indices]
74     class_sample_count = np.array([len(np.where(y_train == t)[0]) for t in np.unique(
        y_train)])
75
76     # Find weights per class
77     weight = 1. / class_sample_count
78     samples_weight = np.array([weight[t] for t in y_train])
79     samples_weight = torch.from_numpy(samples_weight)
80
81     # Define sampler
82     sampler = WeightedRandomSampler(samples_weight.type('torch.DoubleTensor'), len(
        samples_weight))
83
```

```
84     return sampler
85
86 # Create a balanced sampler
87 sampler = balanced_sampler(dataset, train_dataset)
88 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, sampler=sampler)
89
90 dataloaders = {"train": train_dataloader,
91               "val": val_dataloader}
92
93 dataset_sizes = {"train": len(train_dataset),
94                 "val": len(test_dataset)}
95
96
97 # Originally taken from the Pytorch tutorial by Sasank Chilamkurthy
98 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
99     train_losses = []
100     val_losses = []
101     since = time.time()
102
103     # Create a temporary directory to save training checkpoints
104     with TemporaryDirectory() as tempdir:
105         best_model_params_path = os.path.join(tempdir, 'best_model_params.pt')
106
107         torch.save(model.state_dict(), best_model_params_path)
108         best_acc = 0.0
109
110         for epoch in range(num_epochs):
111             print(f'Epoch {epoch}/{num_epochs - 1}')
112             print('-' * 10)
113
114             # Each epoch has a training and validation phase
115             for phase in ['train', 'val']:
116                 if phase == 'train':
117                     model.train() # Set model to training mode
118                 else:
119                     model.eval() # Set model to evaluate mode
120
121                 running_loss = 0.0
122                 running_corrects = 0
123
124                 # Iterate over data.
125                 for inputs, labels in dataloaders[phase]:
126                     inputs = inputs.to(device)
127                     labels = labels.to(device)
128
129                     # zero the parameter gradients
```



```
130         optimizer.zero_grad()
131
132         # forward
133         # track history if only in train
134         with torch.set_grad_enabled(phase == 'train'):
135             outputs = model(inputs)
136             _, preds = torch.max(outputs, 1)
137             loss = criterion(outputs, labels)
138
139             # backward + optimize only if in training phase
140             if phase == 'train':
141                 loss.backward()
142                 optimizer.step()
143
144             # statistics
145             running_loss += loss.item() * inputs.size(0)
146             running_corrects += torch.sum(preds == labels.data)
147         if phase == 'train':
148             scheduler.step()
149
150     epoch_loss = running_loss / dataset_sizes[phase]
151     epoch_acc = running_corrects.double() / dataset_sizes[phase]
152
153     if phase == 'train':
154         train_losses.append(epoch_loss)
155     else:
156         val_losses.append(epoch_loss)
157     print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
158
159     # deep copy the model
160     if phase == 'val' and epoch_acc > best_acc:
161         best_acc = epoch_acc
162         torch.save(model.state_dict(), best_model_params_path)
163
164     print()
165
166     time_elapsed = time.time() - since
167     print(f'Training complete in {time_elapsed // 60:.0f}m {time_elapsed % 60:.0f}
168 }s\n\n')
169     print(f'The best val accuracy score is: {best_acc:.4f}\n\n')
170
171     # Plot the losses
172     plt.plot(train_losses, 'b', label='Training Loss')
173     plt.plot(val_losses, 'r—', label='Validation Loss')
174     plt.legend(loc='upper right')
175     plt.xlabel('Epoch')
```

```
175     plt.ylabel('Loss')
176     plt.title("Model Loss over Epochs for LR=" + str(lr))
177     plt.savefig("losses_"+str(lr)+".png", bbox_inches='tight')
178     plt.close()
179
180     # load best model weights
181     model.load_state_dict(torch.load(best_model_params_path))
182     return model
183
184
185 # We load Resnet since this had the best base accuracy
186 model_conv = torchvision.models.resnet18(weights='IMAGENET1K_V1')
187
188 # This part does the training on the final layer only
189 for param in model_conv.parameters():
190     param.requires_grad = False
191
192 # Parameters of newly constructed modules have requires_grad=True by default
193 num_fts = model_conv.fc.in_features
194 model_conv.fc = nn.Linear(num_fts, 2)
195
196 model_conv = model_conv.to(device)
197
198 criterion = nn.CrossEntropyLoss()
199
200 # Here we try some different learning rates
201 learning_rates = [0.0005, 0.001, 0.005, 0.01, 0.05, 0.1]
202
203 for learning_rate in learning_rates:
204     lr = learning_rate
205     optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=lr, momentum=0.9)
206
207     # Decay LR by a factor of 0.1 every 7 epochs
208     exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
209
210     # Now train the model, and see the loss and accuracy scores for the different
211     # learning rates
212     model_conv = train_model(model_conv, criterion, optimizer_conv,
                              exp_lr_scheduler, num_epochs=num_epochs)
```

Code Listing A.7: code/b\_CNN\_model\_ResNet\_wrs\_momentum.py

```
1 # Import libraries
2 import torch
3 import torch.nn as nn
4 import torch.optim as optim
5 from torch.optim import lr_scheduler
6 import torch.backends.cudnn as cudnn
7 import numpy as np
8 import torchvision
9 from torchvision import datasets, models, transforms
10 import matplotlib.pyplot as plt
11 import time
12 import os
13 from PIL import Image
14 from tempfile import TemporaryDirectory
15 from torch.utils.data import DataLoader
16 import pandas as pd
17 from torch.utils.data import WeightedRandomSampler
18 from sklearn import metrics
19
20 cudnn.benchmark = True
21 plt.ion() # interactive mode
22
23 # Set 'random' seed
24 torch.manual_seed(220029955)
25
26 # Welcome message
27 print("Welcome! We will train the last layer of a pre-trained CNN model.\n")
28
29 # Define the transforms needed
30 data_transforms = transforms.Compose([
31     transforms.Resize([224,224]), # Minimum size needed for Densenet
32     transforms.ToTensor(),
33     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
34     normalisation for Densenet
35 ])
36
37 # Get the dataset from the images created from the wav files
38 dataset = datasets.ImageFolder(os.path.join("data", "images"), transform=
39     data_transforms)
40
41 # Define the classes (Insular and Pelagic)
42 classes = dataset.classes
```

```
43 train_size = int(0.6 * len(dataset))
44 val_size = int((len(dataset) - train_size) / 2)
45 test_size = val_size
46 train_dataset, val_dataset, test_dataset = torch.utils.data.random_split(dataset, [
    train_size, val_size, test_size])
47 print(f"The dataset consists of {train_size + val_size + test_size} datapoints, split
    as follows:")
48 print(f"Train set: {train_size} \nValidation set: {val_size} \nTest size: {test_size}
    \n")
49
50 # Define the device to be used for training
51 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
52 print(f'The device being used is: {device}\n')
53
54 # Define the batch size and number of epochs based on the device
55 if str(device) == "cuda:0":
56     batch_size = 64
57     num_epochs = 24
58 else:
59     batch_size = 20
60     num_epochs = 3
61
62 # Dataloaders
63 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
64 val_dataloader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
65 test_dataloader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
66
67 # Define function to create a balanced sampler
68 # https://www.maskaravivek.com/post/pytorch-weighted-random-sampler/
69 def balanced_sampler(full_dataset, train_dataset):
70     # Find number of samples per class
71     y_train_indices = train_dataset.indices
72     y_train = [full_dataset.targets[i] for i in y_train_indices]
73     class_sample_count = np.array([len(np.where(y_train == t)[0]) for t in np.unique(
        y_train)])
74
75     # Find weights per class
76     weight = 1. / class_sample_count
77     samples_weight = np.array([weight[t] for t in y_train])
78     samples_weight = torch.from_numpy(samples_weight)
79
80     # Define sampler
81     sampler = WeightedRandomSampler(samples_weight.type('torch.DoubleTensor'), len(
        samples_weight))
82
83     return sampler
```

```
84
85 # Create a balanced sampler
86 sampler = balanced_sampler(dataset, train_dataset)
87 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, sampler=sampler)
88
89 dataloaders = {"train": train_dataloader,
90               "val": val_dataloader}
91
92 dataset_sizes = {"train": len(train_dataset),
93                 "val": len(test_dataset)}
94
95
96 # Originally taken from the Pytorch tutorial by Sasank Chilamkurthy
97 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
98     train_losses = []
99     val_losses = []
100     since = time.time()
101
102     # Create a temporary directory to save training checkpoints
103     with TemporaryDirectory() as tempdir:
104         best_model_params_path = os.path.join(tempdir, 'best_model_params.pt')
105
106         torch.save(model.state_dict(), best_model_params_path)
107         best_acc = 0.0
108
109         for epoch in range(num_epochs):
110             print(f'Epoch {epoch}/{num_epochs - 1}')
111             print('-' * 10)
112
113             # Each epoch has a training and validation phase
114             for phase in ['train', 'val']:
115                 if phase == 'train':
116                     model.train() # Set model to training mode
117                 else:
118                     model.eval() # Set model to evaluate mode
119
120                 running_loss = 0.0
121                 running_corrects = 0
122
123                 # Iterate over data.
124                 for inputs, labels in dataloaders[phase]:
125                     inputs = inputs.to(device)
126                     labels = labels.to(device)
127
128                     # zero the parameter gradients
129                     optimizer.zero_grad()
```

```
130
131     # forward
132     # track history if only in train
133     with torch.set_grad_enabled(phase == 'train'):
134         outputs = model(inputs)
135         _, preds = torch.max(outputs, 1)
136         loss = criterion(outputs, labels)
137
138     # backward + optimize only if in training phase
139     if phase == 'train':
140         loss.backward()
141         optimizer.step()
142
143     # statistics
144     running_loss += loss.item() * inputs.size(0)
145     running_corrects += torch.sum(preds == labels.data)
146 if phase == 'train':
147     scheduler.step()
148
149 epoch_loss = running_loss / dataset_sizes[phase]
150 epoch_acc = running_corrects.double() / dataset_sizes[phase]
151
152 if phase == 'train':
153     train_losses.append(epoch_loss)
154 else:
155     val_losses.append(epoch_loss)
156 print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
157
158 # deep copy the model
159 if phase == 'val' and epoch_acc > best_acc:
160     best_acc = epoch_acc
161     torch.save(model.state_dict(), best_model_params_path)
162
163 print()
164
165 time_elapsed = time.time() - since
166 print(f'Training complete in {time_elapsed // 60:.0f}m {time_elapsed % 60:.0f}
167 s\n\n')
168 print(f'The best val accuracy score is: {best_acc:.4f}\n\n')
169
170 # Plot the losses
171 plt.plot(train_losses, 'b', label='Training Loss')
172 plt.plot(val_losses, 'r--', label='Validation Loss')
173 plt.legend(loc='upper right')
174 plt.xlabel('Epoch')
175 plt.ylabel('Loss')
```

```
175     plt.title("Model Loss over Epochs for momentum=" + str(momentum))
176     plt.savefig("losses_"+str(momentum)+".png", bbox_inches='tight')
177     plt.close()
178
179     # load best model weights
180     model.load_state_dict(torch.load(best_model_params_path))
181     return model
182
183
184 # Load the model
185 model_conv = torchvision.models.resnet18(weights='IMAGENET1K_V1')
186
187 # This part does the training on the final layer only
188 for param in model_conv.parameters():
189     param.requires_grad = False
190
191 # Parameters of newly constructed modules have requires_grad=True by default
192 num_fts = model_conv.fc.in_features
193 model_conv.fc = nn.Linear(num_fts, 2)
194
195 model_conv = model_conv.to(device)
196
197 criterion = nn.CrossEntropyLoss()
198
199 # Here we try different momentum values with the lr set to the best lr
200 momentum_values = [0.90, 0.95, 0.97, 0.99]
201
202 for value in momentum_values:
203     momentum = value
204     optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=0.01, momentum=momentum
205                                )
206
207     # Decay LR by a factor of 0.1 every 7 epochs
208     exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
209
210     # Now train the model and see the loss and accuracy scores
211     model_conv = train_model(model_conv, criterion, optimizer_conv,
212                              exp_lr_scheduler, num_epochs=num_epochs)
213
214 # This part creates the classes and scores from the validation data for the metrics
215 y_score = []
216 true_classes = []
217 predicted_classes = []
218
219 for inputs, labels in val_dataloader:
220     inputs = inputs.to(device)
```

```
220     labels = labels.to(device)
221
222     # This part 'flattens' the tensor into a list
223     labels = labels.cpu().numpy().tolist()
224     true_classes.extend(labels)
225
226     with torch.no_grad():
227         model_conv.eval()
228         output = model_conv(inputs)
229
230         for each_output in output:
231             predicted_class = each_output.cpu().data.numpy().argmax() # Numpify each
232             # output
233             predicted_classes.append(predicted_class) # Conactenate
234
235             p = torch.nn.functional.softmax(output, dim=1) # Get probabilities
236             top_proba = p.cpu().numpy()[0][0] # Get probabilities of the positive
237             # class ('Insular') only
238             y_score.append(top_proba)# predicted_classes
239
240 # Let's find the precision, recall and confusion matrix as well
241 print(f"Precision: {round(metrics.precision_score(true_classes, predicted_classes),5)}")
242 print(f"Recall: {round(metrics.recall_score(true_classes, predicted_classes),5)}")
243 print("Confusion matrix:\n", metrics.confusion_matrix(true_classes, predicted_classes))
```



Code Listing A.8: code/b\_CNN\_model\_ResNet\_best.py

```
1 # Import libraries
2 import torch
3 import torch.nn as nn
4 import torch.optim as optim
5 from torch.optim import lr_scheduler
6 import torch.backends.cudnn as cudnn
7 import numpy as np
8 import torchvision
9 from torchvision import datasets, models, transforms
10 import matplotlib.pyplot as plt
11 import time
12 import os
13 from PIL import Image
14 from tempfile import TemporaryDirectory
15 from sklearn import metrics
16 from torch.utils.data import DataLoader
17 import pandas as pd
18 from torch.utils.data import WeightedRandomSampler
19
20 cudnn.benchmark = True
21 plt.ion() # interactive mode
22
23 # Set 'random' seed
24 torch.manual_seed(220029955)
25
26 # Welcome message
27 print("Welcome! We will train the last layer of a pre-trained CNN model.\n")
28
29 # Define the transforms needed
30 data_transforms = transforms.Compose([
31     transforms.Resize([224,224]), # Minimum size needed for Densenet
32     transforms.ToTensor(),
33     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
34     normalisation for Densenet
35 ])
36
37 # Get the dataset from the images created from the wav files
38 dataset = datasets.ImageFolder(os.path.join("data", "images"), transform=
39     data_transforms)
40
41 # Define the classes (Insular and Pelagic)
42 classes = dataset.classes
```

```
43 train_size = int(0.6 * len(dataset))
44 val_size = int((len(dataset) - train_size) / 2)
45 test_size = val_size
46 train_dataset, val_dataset, test_dataset = torch.utils.data.random_split(dataset, [
    train_size, val_size, test_size])
47 print(f"The dataset consists of {train_size + val_size + test_size} datapoints, split
    as follows:")
48 print(f"Train set: {train_size} \nValidation set: {val_size} \nTest size: {test_size}
    \n")
49
50 # Define the device to be used for training
51 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
52 print(f'The device being used is: {device}\n')
53
54 # Define the batch size and number of epochs based on the device
55 if str(device) == "cuda:0":
56     batch_size = 64
57     num_epochs = 24
58 else:
59     batch_size = 20
60     num_epochs = 3
61
62 # Dataloaders
63 val_dataloader = DataLoader(val_dataset, batch_size=batch_size, shuffle=True)
64 test_dataloader = DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
65
66 # Define function to create a balanced sampler
67 # https://www.maskaravivek.com/post/pytorch-weighted-random-sampler/
68 def balanced_sampler(full_dataset, train_dataset):
69     # Find number of samples per class
70     y_train_indices = train_dataset.indices
71     y_train = [full_dataset.targets[i] for i in y_train_indices]
72     class_sample_count = np.array([len(np.where(y_train == t)[0]) for t in np.unique(
        y_train)])
73
74     # Find weights per class
75     weight = 1. / class_sample_count
76     samples_weight = np.array([weight[t] for t in y_train])
77     samples_weight = torch.from_numpy(samples_weight)
78
79     # Define sampler
80     sampler = WeightedRandomSampler(samples_weight.type('torch.DoubleTensor'), len(
        samples_weight))
81
82     return sampler
83
```

```
84 # Create a balanced sampler
85 sampler = balanced_sampler(dataset, train_dataset)
86 train_dataloader = DataLoader(train_dataset, batch_size=batch_size, sampler=sampler)
87
88 ## Visualise the balanced dataset per batch
89 # plot_classes(train_dataloader)
90
91 dataloaders = {"train": train_dataloader,
92               "val": val_dataloader}
93
94 dataset_sizes = {"train": len(train_dataset),
95                 "val": len(test_dataset)}
96
97
98 # Originally taken from the {ytorch tutorial by Sasank Chilamkurthy
99 def train_model(model, criterion, optimizer, scheduler, num_epochs=25):
100     train_losses = []
101     val_losses = []
102     since = time.time()
103
104     # Create a temporary directory to save training checkpoints
105     with TemporaryDirectory() as tempdir:
106         best_model_params_path = os.path.join(tempdir, 'best_model_params.pt')
107
108         torch.save(model.state_dict(), best_model_params_path)
109         best_acc = 0.0
110
111         for epoch in range(num_epochs):
112             print(f'Epoch {epoch}/{num_epochs - 1}')
113             print('-' * 10)
114
115             # Each epoch has a training and validation phase
116             for phase in ['train', 'val']:
117                 if phase == 'train':
118                     model.train() # Set model to training mode
119                 else:
120                     model.eval() # Set model to evaluate mode
121
122                 running_loss = 0.0
123                 running_corrects = 0
124
125                 # Iterate over data.
126                 for inputs, labels in dataloaders[phase]:
127                     inputs = inputs.to(device)
128                     labels = labels.to(device)
```

```
130         # zero the parameter gradients
131         optimizer.zero_grad()
132
133         # forward
134         # track history if only in train
135         with torch.set_grad_enabled(phase == 'train'):
136             outputs = model(inputs)
137             _, preds = torch.max(outputs, 1)
138             loss = criterion(outputs, labels)
139
140         # backward + optimize only if in training phase
141         if phase == 'train':
142             loss.backward()
143             optimizer.step()
144
145         # statistics
146         running_loss += loss.item() * inputs.size(0)
147         running_corrects += torch.sum(preds == labels.data)
148     if phase == 'train':
149         scheduler.step()
150
151     epoch_loss = running_loss / dataset_sizes[phase]
152     epoch_acc = running_corrects.double() / dataset_sizes[phase]
153
154     if phase == 'train':
155         train_losses.append(epoch_loss)
156     else:
157         val_losses.append(epoch_loss)
158     print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
159
160
161
162     # deep copy the model
163     if phase == 'val' and epoch_acc > best_acc:
164         best_acc = epoch_acc
165         torch.save(model.state_dict(), best_model_params_path)
166
167     print()
168
169     time_elapsed = time.time() - since
170     print(f'Training complete in {time_elapsed // 60:.0f}m {time_elapsed % 60:.0f}
171     s\n\n')
172     print(f'The best val accuracy score is: {best_acc:.4f}\n\n')
173
174     # Plot the losses
175     plt.plot(train_losses, 'b', label='Training Loss')
```

```
175     plt.plot(val_losses, 'r—', label='Validation Loss')
176     plt.legend(loc='upper right')
177     plt.xlabel('Epoch')
178     plt.ylabel('Loss')
179     plt.title("Model Loss over Epochs for LR=" + str(lr))
180     plt.savefig("losses_" + str(lr) + ".png", bbox_inches='tight')
181     plt.close()
182
183     # load best model weights
184     model.load_state_dict(torch.load(best_model_params_path))
185     return model
186
187
188 # Load the model
189 model_conv = torchvision.models.resnet18(weights='IMAGENET1K_V1')
190
191 # This part does the training on the final layer only
192 for param in model_conv.parameters():
193     param.requires_grad = False
194
195 # Parameters of newly constructed modules have requires_grad=True by default
196 num_fts = model_conv.fc.in_features
197 model_conv.fc = nn.Linear(num_fts, 2)
198
199 model_conv = model_conv.to(device)
200
201 criterion = nn.CrossEntropyLoss()
202
203 # Set the learning rate and momentum to the best ones found during hyperparameter
    tuning
204 lr = 0.01
205 momentum = 0.95
206 optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=lr, momentum=momentum)
207
208 # Decay LR by a factor of 0.1 every 7 epochs
209 exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
210
211 # Now train the model!
212 model_conv = train_model(model_conv, criterion, optimizer_conv,
213                           exp_lr_scheduler, num_epochs=num_epochs)
214
215 # This part saves entire model
216 torch.save(model_conv, "best-model-resnet.pth")
217
218 # This part tests the model
219 y_score = []
```

```
220 true_classes = []
221 predicted_classes = []
222
223 for inputs, labels in test_dataloader:
224     inputs = inputs.to(device)
225     labels = labels.to(device)
226
227     # This part 'flattens' the tensor into a list
228     labels = labels.cpu().numpy().tolist()
229     true_classes.extend(labels)
230
231     with torch.no_grad():
232         model_conv.eval()
233         output = model_conv(inputs)
234
235         for each_output in output:
236             predicted_class = each_output.cpu().data.numpy().argmax() # Numpify each
237             # output
238             predicted_classes.append(predicted_class) # Conactenate
239
240             p = torch.nn.functional.softmax(output, dim=1) # Get probabilities
241             top_proba = p.cpu().numpy()[0][0]
242             y_score.append(top_proba)# predicted_classes
243
244 # Now let's check the metrics
245 print("\nNow we'll test using the unseen test set. The following are the results:")
246 print(f"Accuracy: {round(metrics.accuracy_score(true_classes, predicted_classes),5)}")
247 print(f"Precision: {round(metrics.precision_score(true_classes, predicted_classes),5)}")
248 print(f"Recall: {round(metrics.recall_score(true_classes, predicted_classes),5)}")
249 print(f"F1: {round(metrics.f1_score(true_classes, predicted_classes),5)}")
250 print("Confusion matrix:\n", metrics.confusion_matrix(true_classes, predicted_classes))
```

Code Listing A.9: code/c\_make\_prediction.py

```
1 # Import Libraries
2 import torch
3 import os.path
4 from scipy import signal
5 from pydub import AudioSegment
6 import numpy as np
7 import matplotlib.pyplot as plt
8 from PIL import Image
9 import sys
10 from torchvision import transforms
11 import warnings
12 warnings.filterwarnings("ignore")
13
14 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
15 print(f'The device being used is: {device}\n')
16
17 # Load the best model
18 if str(device) == "cuda:0":
19     model = torch.load('best-model-resnet.pth')
20 else:
21     model = torch.load('best-model-resnet.pth', map_location=torch.device('cpu'))
22 model.eval()
23
24 # Define function for converting single wav to png file
25 def new_audio_to_png (wavfile):
26     sound = AudioSegment.from_wav(wavfile)
27     sound = sound.set_channels(1)
28
29     # Convert the result to ndarray, and find the sample_rate
30     samples = sound.get_array_of_samples()
31     samples = np.array(samples).astype(np.int16)
32     sample_rate = sound.frame_rate
33
34     # Convert to spectrogram data
35     frequencies, times, spectrogram = signal.spectrogram(samples, sample_rate)
36
37     # Plot the result without any axes or labels
38     plt.pcolormesh(times, frequencies, np.log(spectrogram))
39     plt.tick_params(axis='both', which='both', bottom=False, left=False, top=False,
40                    labelbottom=False, labelleft=False)
41     plt.rcParams["figure.figsize"] = (12,8)
42     plt.savefig("spectrogram_to_be_predicted.png", bbox_inches='tight')
43     plt.close()
44
```

```
45 file_name = sys.argv[1]
46 new_audio_to_png(os.path.join(file_name))
47
48 # Define the transforms needed
49 data_transforms = transforms.Compose([
50     transforms.Resize([224,224]), # Minimum size needed for Resnet
51     transforms.ToTensor(),
52     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225]) # Required
53     normalisation for Resnet
54 ])
55
56 # Define classes
57 classes = [ 'Insular', 'Pelagic' ]
58
59 # Convert the image to a tensor
60 img = Image.open("spectrogram_to_be_predicted.png").convert('RGB')
61 new_tensor = data_transforms(img)
62 new_tensor = new_tensor.unsqueeze(0)
63
64 # Move the data to the GPU if using gpu
65 if str(device) == "cuda:0":
66     new_tensor = new_tensor.to(device)
67
68 # Get prediction
69 with torch.no_grad():
70     output = model(new_tensor)
71     index = output.data.cpu().numpy().argmax()
72     class_name = classes[index]
73
74 # Get probabilities of the prediction
75 p = torch.nn.functional.softmax(output, dim=1)
76 top_proba = p.cpu().numpy()[0][index]
77
78 # Print the results
79 print("output: ", output)
80 print("index:", index)
81 print("class name:", class_name)
82 print("probabilities:", p)
83 print(f"\nThere's a {round((top_proba*100),2)} % chance that this audio file belongs
84     to the {class_name} group.")
```