

Happywhale - Whale and Dolphin Identification

Matthew Harding

April 2024

1 Introduction

This project explores the problem of classifying images of dolphin and whales to identify individuals. This work is based on the *Happywhale - Whale and Dolphin Identification* Kaggle competition [2].

Data for this competition contains over 50,000 images of over 15,000 unique individual marine mammals from 30 different species collected from 28 different research organizations. Individuals have been manually identified and given an `individual_id` by marine researchers.

Unlike a typical classification problem in which there is a fixed set of classes which examples of all classes within the training data, this problem required the ability to classify individuals not contained within the training dataset.

2 Training Data

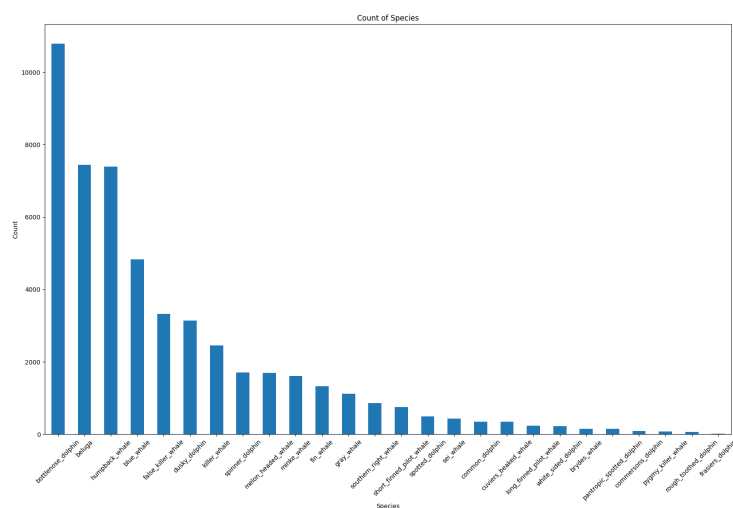


Figure 1: Number of training images per species

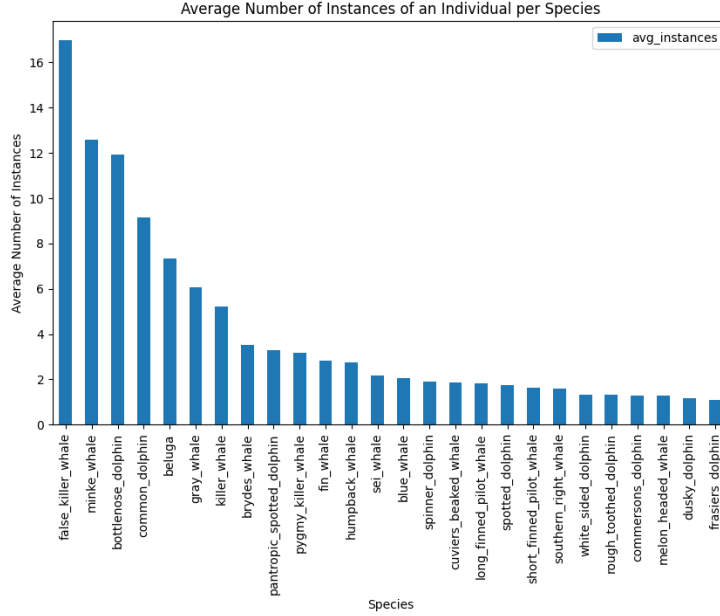


Figure 2: Mean count of images per individual by species

Looking at the histogram of species Figure 1, there is a clear imbalance in the number of images per species with the vast majority of labelled images being for bottlenose dolphins, beluga whales, humpback whales and blue whales. For some species the number of examples is so low it will be difficult to gain a high level of accuracy within classification.

Figure 2 shows the mean count of images per individual by species where we see another imbalance with some species on average containing a large number of images per individual whilst others only contain one or two images per individual.

3 Image Preprocessing

The training images came in a variety of sizes. Figure 3 shows an example unprocessed image. Large amounts of each image are just of water which is irrelevant to the classification task. We are only interested in the parts of the images that contain the animals.

A solution to this issue was proposed within the discussions of the Kaggle competition [1]. YoloV5 was trained on a dataset of 1200 pictures of whale flukes and the corresponding location of points on the edge of the fluke for those pictures [4]. This was then assumed to generalise to finding the bounding box for any part of a whale or dolphin within the water. However, it is worth noting that finding the bounding box of any part of the animal is considered an Out of Distribution problem. From this technique a cropped dataset was



Figure 3: Example image of a bottlenose dolphin

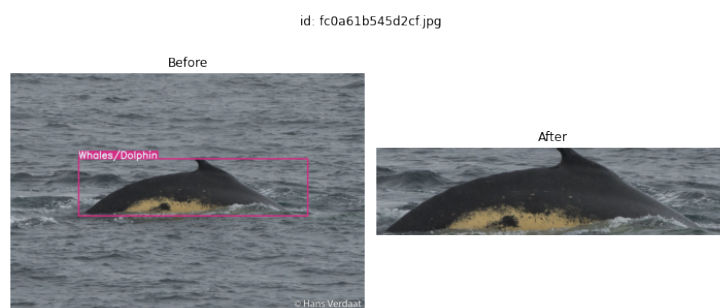


Figure 4: Animal detection using YoloV5 model

created with all images having a dimension of 256 x 256. This is the dataset that was then used to train the whale and dolphin classifier.

4 Image Transformation

A custom PyTorch dataset class called *WhaleDataset* was created for loading in the training images. One-hot encoding was used for the image labels, representing the individual ids provided.

The images were converted to RGB format as this was the compatible format for PyTorch. The images were then transformed using the TorchVision transforms package. The images were already all the same dimension from the preprocessing outlined in the previous section. The transformation involved converting the PIL image to a PyTorch tensor, scaling the pixel values to the range [0,1] followed by normalizing the tensor by subtracting the mean and dividing by the standard deviation.

Scikit Learn's *train_test_split* was used to split of 80/20 was used. A random state was specified for reproducibility between training runs. Shuffling was enabled on the test set.

5 Measuring Model Performance

Performance of the model was measured via the Mean Average Precision @ 5 (MAP@5). Instead of using just accuracy as a measure, MAP@5 was used to account of the complexity of the classification task. Accuracy only considers whether the top predicted class is correct or not. It does not take into account the ranking of predictions. MAP@5, on the other hand, considers the ranking of predictions within the top 5 predicted classes. It rewards models that rank the correct classes higher in the prediction list.

$$MAP@5 = \frac{1}{U} \sum_{u=1}^U \sum_{k=1}^{\min(n,5)} P(k) \times \text{rel}(k)$$

- U is the number of images
- P(k) is the precision at cutoff k
- n is an indicator function equaling 1 if the item at rank k is the correct label, zero otherwise

6 Network Architecture

A convolutional network architecture was selected for this problem. Three convolutional layers were selected for the network architecture so that increasingly more complex and abstract features could be extracted from the input image.

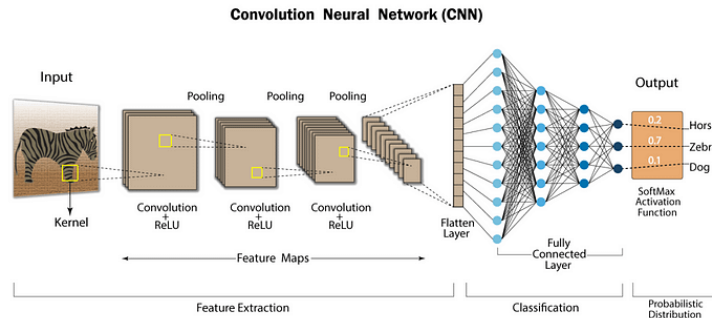


Figure 5: Example of a CNN

- **conv1** takes the input image with 3 color channels (RGB) and applies 32 different 3x3 convolutional filters to extract low-level features like edges, corners, and simple patterns
- **conv2** takes the output of the previous layer (32 feature maps) and applies 64 different 3x3 convolutional filters to detect slightly more complex patterns and features.
- **conv3** takes the output of the previous layer (64 feature maps) and applies 128 different 3x3 convolutional filters to detect even more complex patterns and higher-level features that are useful for distinguishing between different classes.

After each convolutional layer, a ReLU activation function is applied to introduce non-linearity to the network

The first layer was a 2D Convolutional layer. This layer takes the input image with 3 channels (RGB) and applies 32 different 3x3 convolutional filters to extract low-level features like edges, corners and simple patterns. The stride of 1 ensures that the filters are applied to every pixel, and the padding of 1 preserves the spatial dimensionality.

Next, a ReLU activation function was run to add non-linearity to the model followed by a pooling layer that downsamples the feature maps.

After the convolutional layers, there are two fully connected layers. The first fully connected layer takes the flattened output of the convolutional layers and applies a linear transformation, mapping to a fixed-sized vector of 512. This layer acts as a non-linear feature extractor, learning a higher-level representation of the convolutional features. A ReLU activation function is applied after this layer to introduce non-linearity. The final fully connected layer applies a linear transformation, mapping the 512-dimensional vector to a vector of size equal to the number of unique individual ids in the training data.

7 Model Training

The Adam optimizer, a robust gradient-based optimization method. Adam allows for faster convergence, compared to stochastic gradient descent, by adapting the learning rate during training [3]

Despite the massive reduction in file size that came from using the cropped dataset instead of the original dataset, training times on the local machine was still prohibitively long. Therefore, training was done using AWS Sagemaker using a ml.p3.8xlarge GPU instance.

8 Conclusions

It has been shown that a

References

- [1] Happywhale: Cropped data prepare yolov5. Website.
- [2] Happywhale - whale and dolphin identification. Website.
- [3] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- [4] Whale fluke location labelled dataset. Website.