

# Whale Identification Model based on Convolutional Neural Network

Ziyao Qiao qiao.z@husky.neu.edu, Yufei Gao gao.yu@husky.neu.edu  
INFO 7390 Section 04 Spring 2018 Northeastern University

## Abstract

In our final project, we applied many different convolutional neural networks to complete a competition held by Kaggle and Happy Whale.<sup>[1]</sup>

We reshaped the dataset with different size based on the model we use, and we have 2 different version: original photo, original photo with highlighted tail. We applied different network on different dataset and cross those methods in some cases. In the end, we found that VGG19 Model train with ‘new-whale’ class have the highest accuracy.

In our final report, we will introduce the background of deep learning as well as the competition at first. Then we will show related works that help us a lot in the second part of the report. After basic introduction, we will show the mechanism and implementation of our model. Finally, we will show the result of the project and draw a conclusion of our research.

## 1. Introduction

Under decades’ development, Convolutional Neural Networks (CNNs)<sup>[2]</sup> have become one of the most popular approaches to solve computer vision problems, including object identification. Recently, this area developed rapidly because of two main reasons.

Thanks to the rapid development of supervised learning, we could deploy various deep learning model on huge amounts of labeled photos. The wide-spread of advanced hardware such as GPUs also helps us to train a large neural network much faster than ever before.

Most of datasets in the competition comes from Happy Whale, a platform which

already uses image process algorithms to category photos collected from public and various survey, and sorted by Kaggle. For now, the dataset contains one training set and one test set, and one csv file for labeling photos in the dataset.

## 2. Related works

Classifying breed of a specific whale is an example of a fine-grained image classification problem. Fine grained classification is challenging due to subtle differences among the images pertaining to each class. Previously, many researchers conduct research in this area, and they brought out many solutions. Those efforts which is being referred in our research project will be briefly discussed below:

### 2.1 Convolutional Neural Network

In 1966, Marvin Minsky at MIT asked his student spend the summer linking a camera to the computer, and figure the disadvantage of Neural Network at that time. Due to the lack of advanced devices and computing ability, this research paused for a long time. However, in 1989, Yann LeCun created one of the very first convolutional neural networks called LeNet<sup>[3]</sup> based on Back propagation theory, which could return feedback for model to improve. This designation had been used to picture recognition early and propelled the field of Deep Learning.

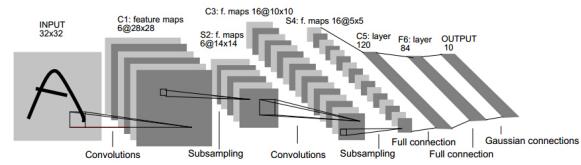


Figure 2.1 Architecture of LeNet

## 2.2 VGG-19<sup>[4]</sup>

In 2014, a team called Visual Geometry Group brought out the VGG Model, which contains deeper network and perform better on image identification problems than other models at that time.

Instead of using so many hyper parameters, all convolution layers in VGG-19 are  $3 \times 3$  filter and 1 strides with same padding.

Overall, VGG19 Model contains 16 convolutional layers and 3 pooling layers, and use “softmax” before return the result. The key that made VGG-19 successful is the depth of the network and the huge numbers of trainable parameters, which is about 138 million.

| ConvNet Configuration  |                        |                               |  |  |  |
|------------------------|------------------------|-------------------------------|--|--|--|
| A                      | A-LRN                  | B                             | C  | D  | E  |
| 11 weight layers       | 11 weight layers       | 13 weight layers              | 16 weight layers                           | 16 weight layers                           | 19 weight layers                           |
| conv3-64               | conv3-64<br>LRN        | conv3-64<br><b>conv3-64</b>   | conv3-64<br>conv3-64                       | conv3-64<br>conv3-64                       | conv3-64<br>conv3-64                       |
| maxpool                |                        |                               |  |  |  |
| conv3-128              | conv3-128              | conv3-128<br><b>conv3-128</b> | conv3-128<br>conv3-128                     | conv3-128<br>conv3-128                     | conv3-128<br>conv3-128                     |
| maxpool                |                        |                               |  |  |  |
| conv3-256<br>conv3-256 | conv3-256<br>conv3-256 | conv3-256<br>conv3-256        | conv3-256<br>conv3-256<br><b>conv1-256</b> | conv3-256<br>conv3-256<br><b>conv1-256</b> | conv3-256<br>conv3-256<br><b>conv1-256</b> |
| maxpool                |                        |                               |  |  |  |
| conv3-512<br>conv3-512 | conv3-512<br>conv3-512 | conv3-512<br>conv3-512        | conv3-512<br>conv3-512<br><b>conv1-512</b> | conv3-512<br>conv3-512<br><b>conv1-512</b> | conv3-512<br>conv3-512<br><b>conv1-512</b> |
| maxpool                |                        |                               |  |  |  |
| FC-4096                | FC-4096                | FC-4096                       | FC-1000                                    |  | soft-max                                   |

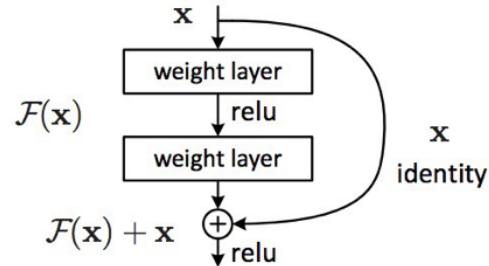
**Figure 2.2 Architecture of VGG Model  
(VGG19 is in the right site of the graph)**

Similar to other VGG Model's features, VGG-19 also follows the pattern that when goes deeper into the network, the number of channel rises while the width and height of the volume reduces.

## 2.3 ResNet<sup>[5]</sup>

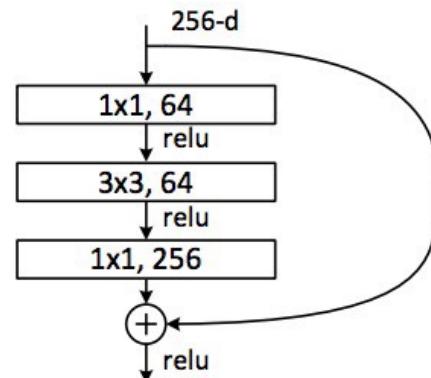
The revolution happened again in December 2015. ResNet have relatively simple ideas: feed the output of two successive

convolutional layer, bypass the input to the next layers, improve and update information from residual value.



**Figure 2.3 Basic Idea of ResNet Model**

This is similar to older ideas like this one. But here they bypass 2 layers and are applied to large scales. Bypassing after 2 layers is a key intuition, as bypassing a single layer did not give much improvements. A 2-layers architecture can be thought as a small classifier. This is also the very first time that a network has more than 1,000 layers was trained.



**Figure 2.4 Residual block in ResNet50 Model**

This layer reduces the number of features at each layer by first using a  $1 \times 1$  convolution with a smaller output (usually  $1/4$  of the input), and then a  $3 \times 3$  layer, and again a  $1 \times 1$  convolution to a larger number of features. Like in the case of Inception modules, this allows to keep the computation low, while providing rich combination of features.

## **2.4 InceptionV3**

In February 2015, close to the same time ResNet created, Inception V2 was introduced to public for first time. Batch-normalization computes the mean and standard-deviation of all feature maps at the output of a layer, and normalizes their responses with these values.

This corresponds to ‘whitening’ the data, and thus making all the neural maps have responses in the same range, and with zero mean. This helps training as the next layer does not have to learn offsets in the input data, and can focus on how to best combine features.

Based on the principle, in December 2015, InceptionV3 had been created, with functions which could catch features more effectively.

## **3. Problem analysis**

### **3.1 Lots of categories with limited training set**

The training-set totally have 4251 categories, including 4250 unique and 1 ambiguous. However, we have only 9851 photos in the dataset, only 2-3 photo per category in average.

### **3.2 Unbalanced dataset distribution**

As we described in midterm presentation, more than 3500 special ID in the dataset have less than 3 photos. And the

“new\_whale”, one ambiguous category, have 810 photos. Those photos also make huge noise, so we would use specific balance methods to limit noise, and we also try remove them in some cases.

## **4. Implement**

### **4.1 Preprocessing**

#### **4.1.1 Computing Environment Set up**

Our laptops have relative low computing ability, so we run our project in a computer with GTX-1080ti GPU and AMD 8-core CPU. We used CUDA-GPU in our training

process, in this way we could accelerate our process and save our time, and run more models in limited time without paying money to Cloud-Computing service provider like Amazon and Google.

#### **4.1.2 Resize photos, recolor and set label**

In the dataset, almost every photo has different size. In addition, some of them are black and white, but others are colored. Thus, we resize them before training. Besides, we also tried to make all the photos black and white before training in some cases.

#### **4.1.3 Highlight Tail Area**

Since tail is only a small part of photos, other objects such as sea-waves, ship and texts could make large noise in the training process. We designed an algorithm to try to highlight the tail area based on features of whales tail<sup>[6]</sup> to improve the accuracy.

#### **4.1.4 Data Augmentation<sup>[7]</sup>**

Since we have limited data in the training set, so we augment data for training and internal testing in our training process. We use different to spin and reverse the photos, and try to figure out if the model catches the features correctly. If it returns an wrong answer, the residual could return some useful information to the model and help it to improve its performance.

## **4.2 Basic Convolutional Model**

In the very beginning we constructed a CNN Model as baseline of the project. We cross validate the model with different elements including architecture, activation function, pool size, loss function, optimizer, batch size, epoch and step. Finally, we set up a model with relative high cost-interest ratio, with 3 convolutional layers, train 128 photos per step, iterate all photos per epoch and train for up to 100 epochs.

#### **4.2.1 Network Architecture**

Initially, our model has 3 convolutional layers, with several dropout layers, flatten layer and dense layer to reshape the data and catch features. We set batch-size as 128 and step value is total number of the training set devided by batchsize, in this way we could make sure every photo in the dataset could be iterated once in a single epoch. Initially, we set the model run in small epochs like 10 to run in relative low costs. However, more epochs could be added if someone have better machine. In the end, we increase the epochs up to 100.

Besides, we change many different methods by using cross validation method. Compare the cost and accuracy, then find out the best one as part of our CNN model<sup>[12]</sup>.

#### **4.2.2 Activation Function**

Various activations functions are used for the convolution layers and the dense layers and the accuracy was calculated. We tried to change the activation function from Rectified Linear Unit (ReLU) and Sigmoid to exponential linear unit (elu) and softplus. However, for this dataset, those two functions seem useless and the accuracy decrease drastically. Thus, we used ReLU and Sigmoid as the activation in our CNN Model.

#### **4.2.3 Cost Function**

Just like activation function, we change the cost function from categorical-crossentropy to mean squared error. The substitution works bad for non-numerical categorical dataset, as it could not provide useful data for model to learn and improve, so the new function doesn't work well in our dataset. Hence, categorical-crossentropy used as the cost function for our model.

#### **4.2.4 Kernel Initializer**

Using various kernel initializers did not produce significant changes in the model

accuracy or performance. Hence, we decided to use the default initializer to keep the performance of the model relatively good.

#### **4.2.5 Optimizer Function**

The model was trained with Adadelta Optimizer. We also tried Adamax and Momentum as optimizer function. The network plateaued normally when we used Adadelta Optimizer but overfitting fastly when we use others. At last we keep using Adadelta in our CNN Model.

#### **4.2.6 Number of Epochs**

Since in every single epoch of our model, all photos in our training set could be iterated for once, so we needn't consider epochs with the quantity of our dataset. At first we just set our epoch as 1 to make sure our model could run normally, however, more epochs were added as we have better machine and in the end, we increase the epochs up to 100.

#### **4.2.7 Learning rate**

The model was trained with 0.001 as the training rate when we had 6 classes. But as the classes increased, the model took a lot of steps to plateau. Hence, we implemented decaying learning rate with the initial learning rate of 0.05 and the decay factor as 0.1, which decays every 50 epochs.

#### **4.3 Pre-trained Model<sup>[10]</sup>**

Keras have official VGG19<sup>[8]</sup>. ResNet50 and InceptionV3 implement<sup>[9]</sup>, so we just applied those pretrained model to our dataset. As what we had done in the basic convolutional model before, we also need to augment data for training and internal testing in our pre-trained model training process.

## 5.Conclusion

### 5.1 Preprocessing Stage

#### 5.1.1 Algorithms concerning detecting contour area

Our program aiming to figure out the tail basically based on “THRESH\_BINARY” and “CANNY” function in cv2 package, two packages which could find angle and figure out the edge of the shape in the picture. However, we underestimated the noise made by sea waves and other contents in the pictures. In some photo which is relatively clean, the algorithm performs well. But in others with various noises, this program seems useless.

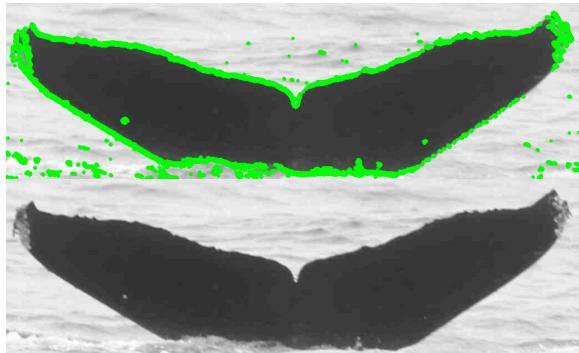


Figure 5.1 A success sample of contour-detect algorithm



Figure 5.2 A failure example because of unexpected angles from spots in tail and waves

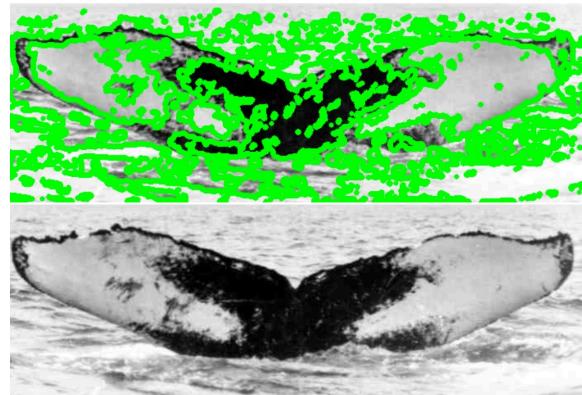


Figure 5.3 Another failure example generated by waves

#### 5.1.2 Computing ability

Besides the model we choose and methods we used to tackle with photos based on our devices, computation power is also an important part of performances in such competitions. Since lack of advanced machines and computing ability of large models, we didn't train those models systematically, and we also unable to try any other aggregated models or tackle with photos with other advanced algorithms, which definitely limited our performance. Furthermore, the average utility of our GPU(GTX-1080ti) always keep 20%-30%, which seldom fully used. As we analyzed, the bottleneck is the GPU Memory. The GPU could calculate a large batch of photo in one time but GPU-Memory limited the data stream<sup>[11]</sup>.

| Method                             | Accuracy in test set | Epochs | Average time (s/per epoch) |
|------------------------------------|----------------------|--------|----------------------------|
| ResNet50 without 'NewWhale'        | 32.482%              | 10     | 12478s                     |
| ResNet50 with 'NewWhale'           | 32.631%              | 10     | 12600s                     |
| VGG19 with 'NewWhale'              | 32.999%              | 30     | 2270s                      |
| InceptionV3 without 'NewWhale'     | 32.481%              | 20     | 4703s                      |
| CNN Model with contour detected    | 32.763%              | 100    | 215s                       |
| CNN Model without contour detected | 32.875%              | 100    | 214s                       |

Figure 5.4 Overall result of different models in our projects

Another point is that we could design some better algorithms in data augmentation. As we see, the test set in training process seems not work very well for the final result, since we just rotate the photos in the training set simply and shift the size, which seems helpless for the model to predict photo from another dataset.

## **5.2 Models in the project and dataset**

In conclusion, the basic CNN model seems too easy for the dataset and it could not figure out too many features from the training set. Besides, since many class in the dataset have only one or two picture, so the weight of features from single photo could be really large and have negative influence to the result.

As for pretrained model, VGG19 seems work best for this project. That's mainly because in the dataset, most of classes have only one or two photos, so a model with too many layers may work worse than a relative simple model. The architecture of VGG19 probably fits the dataset and could catch features better than other two models.

However, even different models we implemented have tiny difference, totally they performed similar for this dataset and the score in the test set not have large difference. As we summarized, that's mainly because of the low quality of the dataset: nearly every pictures have different size, some photos is colorful but others not, and the photos obtained from different direction, different distance and environment, which is also consider and captured as a part of features by the models.

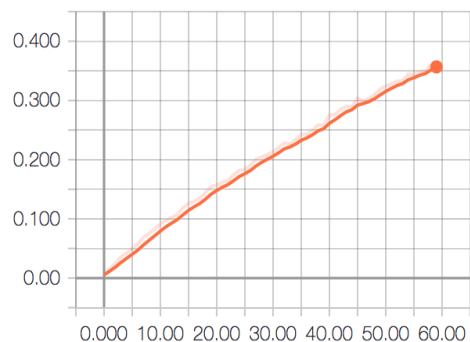
All implementation concerning the project could be found on our GitHub repo.

## Reference

- [1] <https://www.kaggle.com/c/whale-categorization-playground>
- [2] <https://www.coursera.org/learn/neural-networks-deep-learning>
- [3] [https://github.com/BVLC/caffe/blob/master/examples/mnist/lenet\\_train\\_test.prototxt](https://github.com/BVLC/caffe/blob/master/examples/mnist/lenet_train_test.prototxt)
- [4] <https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
- [5] <https://arxiv.org/abs/1512.03385>
- [6] [https://en.wikipedia.org/wiki/Humpback\\_whale](https://en.wikipedia.org/wiki/Humpback_whale)
- [7] <https://cs231n.github.io>
- [8] <https://medium.com/@14prakash/transfer-learning-using-keras-d804b2e04ef8>
- [9] <https://www.learnopencv.com/keras-tutorial-transfer-learning-using-pre-trained-models/>
- [10] [https://www.tensorflow.org/tutorials/image\\_retraining](https://www.tensorflow.org/tutorials/image_retraining)
- [11] <https://github.com/spejss/Keras-Transfer-Learning-Tutorial>
- [12] <https://www.kaggle.com/gimunu/data-augmentation-with-keras-into-cnn>

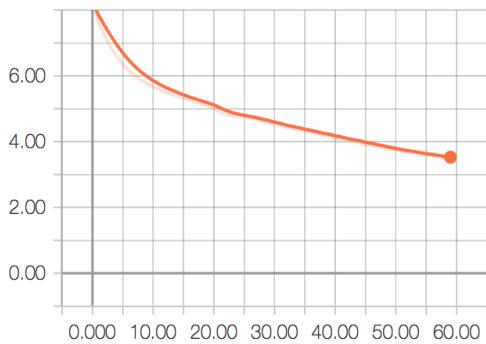
## APPENDIX

acc



loss

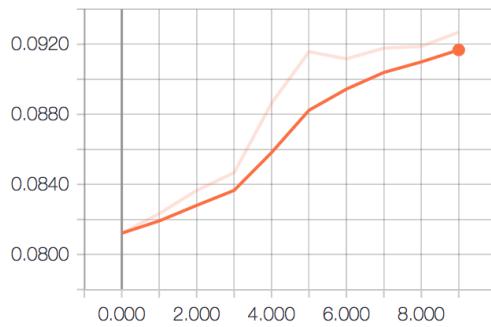
loss



**performance of InceptionV3 in training set in one epoch**

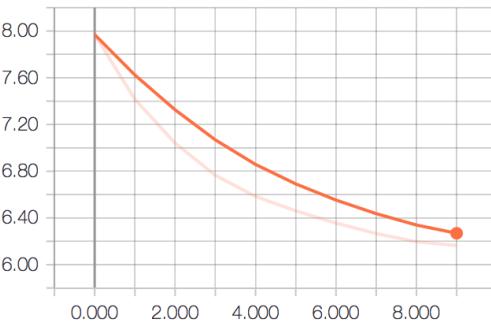
acc

acc



loss

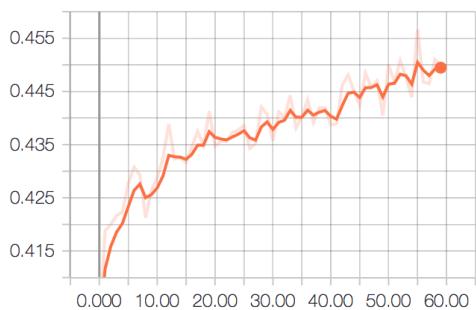
loss



**performance of VGG19 in training set in one epoch**

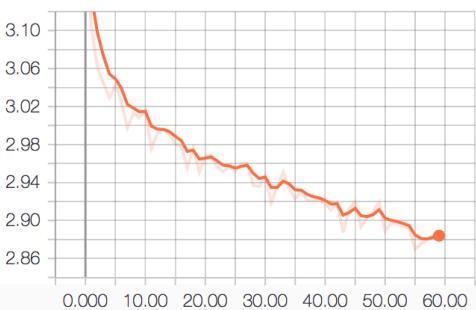
acc

acc



loss

loss



**performance of ResNet50 in training set  
in one epoch**