

Udacity Machine Learning Nanodegree

Capstone Project: Dog Breed Classifier

George Xian Wee

Project Definition

Overview

This project originates from one of the capstone proposals provided by Udacity in the Machine Learning Engineer Nanodegree. It is a dog breed classifier which integrates several components to create a simple application to identify dogs and humans and to classify the type of dog breed, or the type of dog breed a human most closely resembles. The project falls under the research domain of image recognition, a field in machine learning where major breakthroughs have been achieved, such as a convolutional neural network (CNN) winning 2012 ImageNet Large Scale Recognition challenge. The training, validation and test dataset in this project was provided by Udacity.

Problem Statement

The project primarily addresses classification problems using supervised learning. Input images are mostly of dogs and human faces but the application can actually accept any image. The first classification problem is to determine if an image contains a dog (the model will classify an image as dog or not dog). If there is a dog in the image, a convolutional neural network will take the dog image as an input and

output a prediction of the breed of the dog. If there is not a dog in the image, the application will detect if there is a human face in the image and if one has been found, the CNN will guess which kind of dog breed the human face most closely resembles. In the case where an image does not contain a dog or human, the application will display a message noting this state.

Metrics

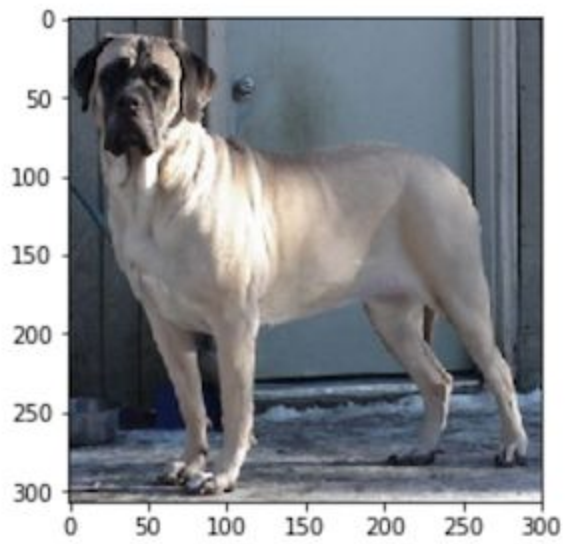
In this project accuracy (number of correct predictions / total number of predictions) is used as the primary metric for evaluation and is suitable because the data is fairly balanced, meaning the training, validation and test data have roughly the same number of samples of dogs of each breed. A balanced dataset is important to avoid misleading accuracy results, such as the case in which 90% of the test data images are of a specific breed - then a default prediction of that breed would yield 90% accuracy.

Analysis

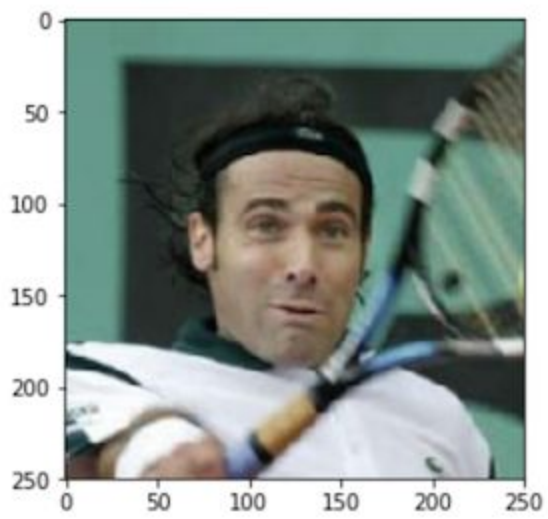
Data Exploration and Visualization

The dataset contains 8351 dog images and 13233 human images. Each dog and human image in its raw form has R G B values and may vary in its pixel dimensions. For example, an image tensor could have the following dimensions: [1, 3, 224, 224], 1 image, a list RGB values of length 3 and width and length of 224.

Sample Dog Image



Sample Human Image



Methodology

Data Preprocessing

The images were retrieved using the Python Image Library (PIL), resized, converted to tensors and normalized. After the image transformations, a sample dog or human image would have, for example, dimensions 1 x 3 x 224 x 224, meaning the first tensor would contain 3 tensors representing RGB values of a 224 x 224 image. Each pixel is an input into the initial layer of the neural network.

For the model built using transfer learning, images were transformed into uniform dimensions by first resizing them to 256 pixels and then center cropping to 224 pixels. Since PyTorch's ResNet model was used as the starting point for transfer learning, I used a similar transformation process and dimensions from its documentation https://pytorch.org/hub/pytorch_vision_resnet/. The images were normalized with mean values [0.485, 0.456, 0.406] and standard deviation values [0.229, 0.224, 0.225], which are also from ResNet's recommended mean and standard deviation values.

For the model built from scratch, the images were resized and cropped to 32 x 32 pixels to reduce training time and normalized with mean values [0.5, 0.5, 0.5] and standard deviation values [0.5, 0.5, 0.5].

Implementation and Refinement

First I defined, trained, validated and tested a convolutional neural network from scratch. The CNN contains 2 layers of convolutions which produce features used by a final, fully connected layer that takes in the features and outputs predictions. The first convolutional layer has 3 input channels to account for the 3 dimensions of R, G, B values, 6 output channels, a kernel size of 5 for each sliding filter and a stride of 1. Generally, odd numbers are chosen for kernel sizes in the convolutional layer to avoid

distortions that occur across layers with even numbered kernel sizes. For example, common kernel sizes in convolutional layers are 3x3, 5x5 and 7x7. Smaller strides tend to encode more information and maintain translational invariance, meaning the effect of the position of the object in the image is reduced. The second convolutional layer has 6 input channels which is the same number of output channels as from layer 1, 16 output channels and a kernel size of 5.

To downsample, or to reduce the dimensionality of each layer, max pooling was used to obtain the maximum value from the activated features in the previous layer (the maximum activation). A kernel size and stride of 2 were chosen, which is consistent with common practice. A small kernel size and stride in pooling prevents discarding too much information from the previous layer.

Finally, fully connected layers take in the features from the convolutional layers and produces predictions. The input layer has $16 * 5 * 5$ parameters (the number of outputs from the last convolutional layer multiplied by the dimension of the filters in the 2nd convolutional layer). The hidden layers consist of 120 and 84 neurons. There are a few empirically driven rules of thumb when choosing the number of hidden neurons. In general, the number of hidden neurons should be between the number of input and output parameters. Too few hidden neurons may result in underfitting and too many hidden neurons may lead to overfitting. A common rule is to have the number of hidden neurons be roughly the average of the number of neurons in the input and output layers. The output layer contains 133 neurons which corresponds to the 133 dog breeds in our dataset. Each neuron of the output layer is a probability assigned to each dog breed, with the highest probability representing the model's prediction of the type of dog breed in the image.

The prediction results of a neural network built from scratch and trained over 25 epochs of 64 images per batch was quite low, only about 2% accuracy (21/836). More details are provided in the Evaluation section below.

To improve the performance, I employed fixed feature extraction, a form of transfer learning to extract the features from a pre-trained neural network, replace the final layer of the network and only optimize the weights of that layer. I started with a pre-trained neural network, ResNet18, a residual neural network, which is a convolutional neural network containing “skip connections” to simplify the network by using fewer layers for training and to avoid “vanishing gradients” problem in which the gradients in the early layers of a neural network become extremely small, inhibiting the network’s ability to accurately reflect how a small change in a parameter’s value will affect the output.

Next I replaced the final prediction layer of the network with a fully connected layer consisting of 512 input features (the output from the final convolutional layer of the model) and 133 output nodes (the number of dog breeds in our dataset). The final layer was re-trained with Cross Entropy Loss criterion and stochastic gradient descent optimizer to update the parameters with a learning rate (how much to adjust each parameter based on the gradient) of .001 and a momentum (the accumulation of previous gradients to determine how much to update the parameters) of .9, a common value used in practice. Learning rates (typically ranging from .1 - .001) that are larger enable models to train faster but may lead to suboptimal final weights while smaller learning rates train more slowly but may lead to better results. The result, discussed in the next section, substantially improved.

Example Transfer Model Training/Validation Loss:

Epoch: 1	Training Loss: 4.672806	Validation Loss: 4.114409
Epoch: 2	Training Loss: 3.725128	Validation Loss: 3.272928
Epoch: 3	Training Loss: 3.020625	Validation Loss: 2.669618
Epoch: 4	Training Loss: 2.502269	Validation Loss: 2.198346
Epoch: 5	Training Loss: 2.132495	Validation Loss: 1.882919
Epoch: 6	Training Loss: 1.849949	Validation Loss: 1.651998
Epoch: 7	Training Loss: 1.644703	Validation Loss: 1.457477

Epoch: 8 Training Loss: 1.480561	Validation Loss: 1.318219
Epoch: 9 Training Loss: 1.350160	Validation Loss: 1.197332
Epoch: 10 Training Loss: 1.242551	Validation Loss: 1.122416
Epoch: 11 Training Loss: 1.159051	Validation Loss: 1.033086
Epoch: 12 Training Loss: 1.080627	Validation Loss: 0.973539
Epoch: 13 Training Loss: 1.016656	Validation Loss: 0.917403
Epoch: 14 Training Loss: 0.966058	Validation Loss: 0.862755
Epoch: 15 Training Loss: 0.916244	Validation Loss: 0.826528
Epoch: 16 Training Loss: 0.876311	Validation Loss: 0.786646
Epoch: 17 Training Loss: 0.843698	Validation Loss: 0.755595
Epoch: 18 Training Loss: 0.800897	Validation Loss: 0.725584
Epoch: 19 Training Loss: 0.778323	Validation Loss: 0.690264
Epoch: 20 Training Loss: 0.746775	Validation Loss: 0.674685
Epoch: 21 Training Loss: 0.726695	Validation Loss: 0.646116
Epoch: 22 Training Loss: 0.700255	Validation Loss: 0.624603
Epoch: 23 Training Loss: 0.678680	Validation Loss: 0.609314
Epoch: 24 Training Loss: 0.660038	Validation Loss: 0.590945
Epoch: 25 Training Loss: 0.638552	Validation Loss: 0.576139

Results

Model Evaluation, Validation and Justification

To evaluate both the model created from scratch and the transfer learning model, I used the following validation process. I accumulated the loss during each forward pass and computed the average validation loss by dividing the total loss in each epoch by the number number of samples in the validation set. The model was validated on 835 images. To track the progress of the training, I accumulated the training loss in each epoch and calculated the average training loss. The model was trained on 6680 images. For the convolutional network built and trained from scratch, the results were poor. After training for 25 epochs, the training loss only decreased from 4.892253 to 4.778899.

Validation loss was reduced from 4.891494 to 4.778537. The results of the model on the test set was similar, a 4.752226 loss. The model was only able to achieve an accuracy of about 2% (21/836).

I then experimented with increasing the learning rate by an order of magnitude, to .01, to see if performance would improve. Over a course of 25 epochs, training loss was reduced to 3.113258 from 4.892251. Validation loss dropped from 4.889027 to 3.119743. Test loss was 0.498703 and the accuracy achieved by a model trained on a higher learning rate improved from 2% to about 8%.

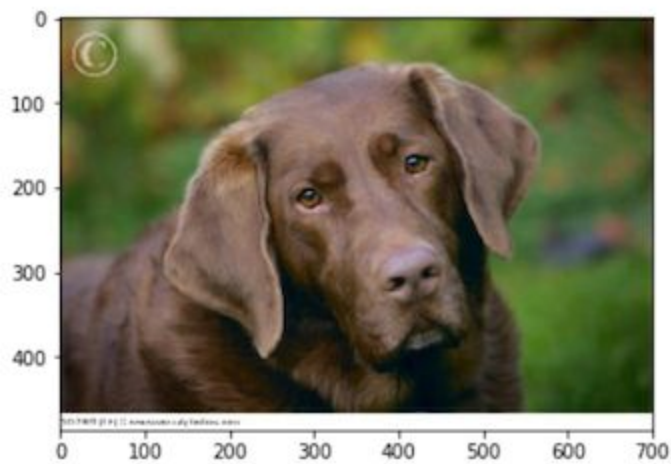
The results from transfer learning were substantially better. In epoch 1, the training loss started at 4.672806 and by epoch 25, was reduced to 0.638552. The validation loss began at 4.114409 and ended at 0.576139. When the trained transfer model was tested on the test dataset, the model's loss was 0.787483 and achieved an accuracy of 82% (684/836). When the learning rate was increased to .01, training loss decreased from 2.898901 to 0.123613, validation loss decreased from 1.440444 to 0.092910 and the model accuracy ticked up to 83%.

lr : learning rate	Test Loss	Accuracy
Model Scratch, lr = .001	4.752226	2%
Model Scratch lr = .01	4.536067	8%
Model Transfer lr = .001	0.787483	82%
Model Transfer lr = .01	0.498703	83%

Example 1: Dog

Ground truth: Labrador Retriever

```
Dog detected!  
Dog breed: Labrador retriever
```



Example 2: Dog

Ground truth, as indicated in file path is Brittany

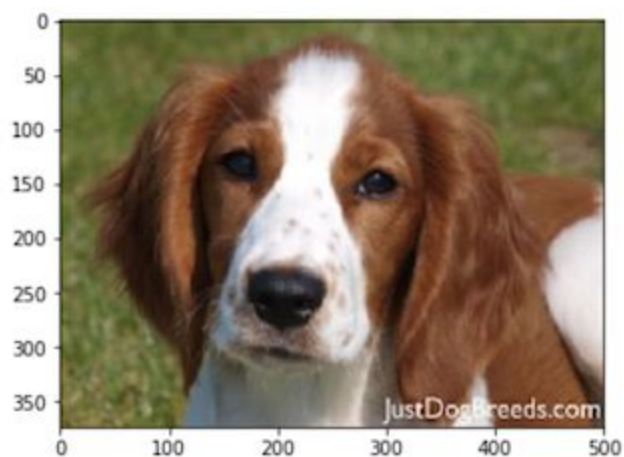
```
Dog detected!  
Dog breed: Brittany
```



Example 3: Dog

Ground truth is Welsh Springer Spaniel

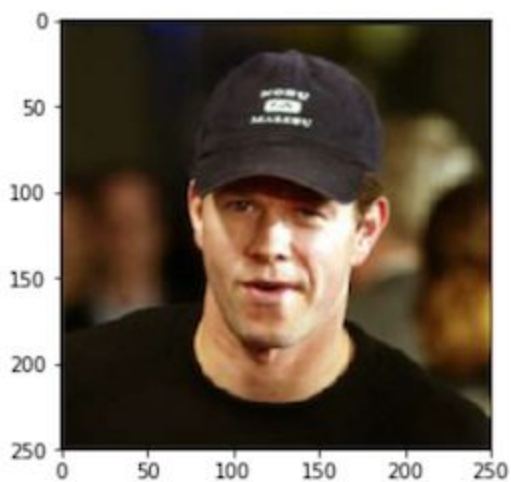
Dog detected!
Dog breed: Welsh springer spaniel



Example 4: Human

Image of actor Mark Wahlberg

Human detected!



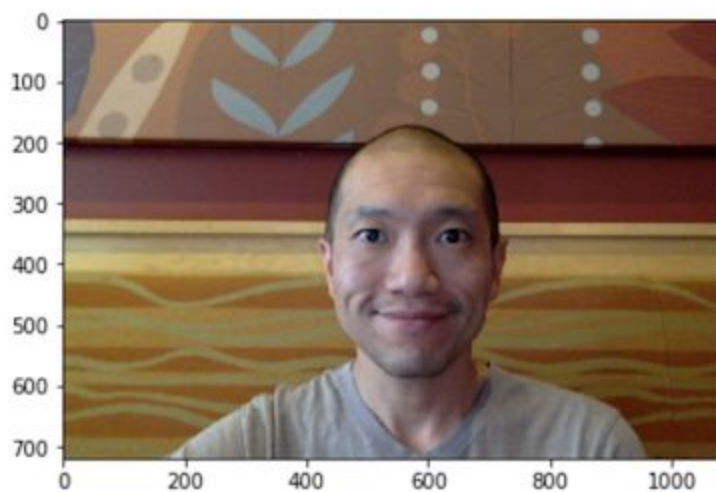
Dog breed resembled: Doberman pinscher



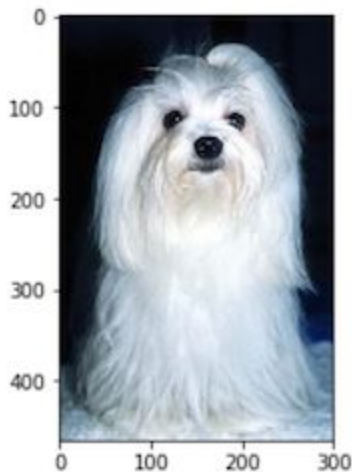
Example 5: Human

Person: George

Human detected!



Dog breed resembled: Maltese



Conclusion

In this project, images of dogs were preprocessed, a convolutional neural network was trained, validated and tested, and a function was written to simulate an application that takes in an input image path, predict the breed of the dog or the type of breed a human most closely resembles. Building a convolutional neural network that can achieve high accuracy from scratch was tougher than I expected. Initially I thought that by increasing the learning rate, I could increase the accuracy, which did occur but only up to a point. To improve the results of the model trained from scratch, I can try adding more convolutional layers and to boost the accuracy of the transfer model, perhaps I can experiment with adding more hidden layers to the final fully connected layers of the network.

Sources:

<http://cs231n.github.io/convolutional-networks/> <http://giant.uji.es/blog/convnet/convnet.html>

<https://stats.stackexchange.com/questions/181/how-to-choose-the-number-of-hidden-layers-and-nodes-in-a-feedforward-neural-net>

<https://towardsdatascience.com/deciding-optimal-filter-size-for-cnns-d6f7b56f9363>

<https://www.quora.com/How-can-I-decide-the-kernel-size-output-maps-and-layers-of-CNN>

<https://arxiv.org/pdf/1606.02228v2.pdf> <https://www.quora.com/How-does-one-determine-stride-size-in-CNN-filters>

<https://stats.stackexchange.com/questions/207195/translational-variance-in-convolutional-neural-networks>

<https://towardsdatascience.com/a-guide-to-an-efficient-way-to-build-neural-network-architectures-part-ii-hyper-parameter-42efca01e5d7>

<https://www.quora.com/Are-maxpooling-layer-kernel-sizes-in-CNNs-generally-smaller-than-convolutional-layer-kernel-sizes-Why>

<https://www.kaggle.com/pytorch/resnet18>

<https://towardsdatascience.com/understanding-and-coding-a-resnet-in-keras-446d7ff84d33>

<https://www.quora.com/What-is-the-vanishing-gradient-problem>

<https://towardsdatascience.com/stochastic-gradient-descent-with-momentum-a84097641a5d>