

CNN Project: Dog Breed Classifier

A Capstone Report

By

PRAVEEN BANDARU

Submitted to Udacity

in partial fulfillment of the requirements for the degree of

Machine Learning Engineer

Nanodegree Program

March 2020

Definition

Project Overview

Image classification [1] is a supervised learning [2] problem: define a set of target classes (objects to identify in images) and train a model to recognize them using labeled example photos. A breakthrough in building models for image classification came with the discovery that a convolutional neural network [3] (CNN) could be used to progressively extract higher- and higher-level representations of the image content. Instead of preprocessing the data to derive features like textures and shapes, a CNN takes just the image's raw pixel data as input and "learns" how to extract these features, and ultimately infer what object they constitute.

In this project, I created a code pipeline using transfer learning to create a CNN that can identify dog breed from images. We used a pretrained ResNet50 model from Torchvision Model Zoo for our Dog Classification problem. This pretrained model was trained on ImageNet dataset and has the last Fully Connected layer with 1000 out features. These pre-trained networks demonstrate a strong ability to generalize to images outside the ImageNet dataset via transfer learning. We make modifications in the pre-existing model by fine-tuning the model. We have frozen the parameters, so we don't back-propagate through them and replaced the last Fully Connected layer with a Linear layer having 133 out features equal to the number of classes in our Dog dataset. Our pipeline will accept any user-supplied image as input. If a dog is detected in the image, it will provide an estimate of the dog's breed. If a human is detected, it will provide an estimate of the dog breed that is most resembling. The image below displays a sample output of dog detection pipeline.

Problem Statement

In this project, we will learn how to build a pipeline to process real-world, user-supplied images. Given an image of a dog, our algorithm will identify an estimate of the canine's breed. If supplied an image of a human, the code will identify the resembling dog breed. Along with exploring state-of-the-art CNN models for classification, we will make important design decisions about the user experience for our app. Our goal is that by

completing this project, we understand the challenges involved in piecing together a series of models designed to perform various tasks in a data processing pipeline. Each model has its strengths and weaknesses, and engineering a real-world application often involves solving many problems without a perfect answer. Our imperfect solution will nonetheless create a fun user experience!

Evaluation Metrics

The model's performance will be evaluated based on the test set accuracy score, i.e., we want to see how the model will perform when it encounters an image it has never seen before. Accuracy is a common metric for image classification, it takes into account both the true positives and true negatives with equal weight.

$$\text{Accuracy} = (\text{true positives} + \text{true negatives}) / \text{dataset size}$$

We achieved a test accuracy of 85% after running our model on the test dataset.

Analysis

Data Exploration

The required human [4] and dog [5] datasets were provided by Udacity as downloadable links in the *dog_app.ipynb* notebook. They can be also found in the /data folder in the Udacity workspace. There are 13233 total human images and 8351 total dog images which spanned across 133 breeds. The dog images were already split into *test*, *train* and *valid* subfolders and each of the *test*, *train* and *valid* folders were further divided into 133 breeds subfolders. In the *dog_app.ipynb* notebook, we save the file paths for both the human (LFW) dataset and dog dataset in the numpy arrays *human_files* and *dog_files*.

Exploratory Visualization

Here is an exploratory visualization of the training data set.

The following plot shows a sample image for the first 50 classes from the training data set.

Sample image in each class



Fig. 1 Images from the first 50 classes of the dog training data set

Next, I plotted a bar chart showing how the data is distributed across different classes. It shows the total count of images from each class.

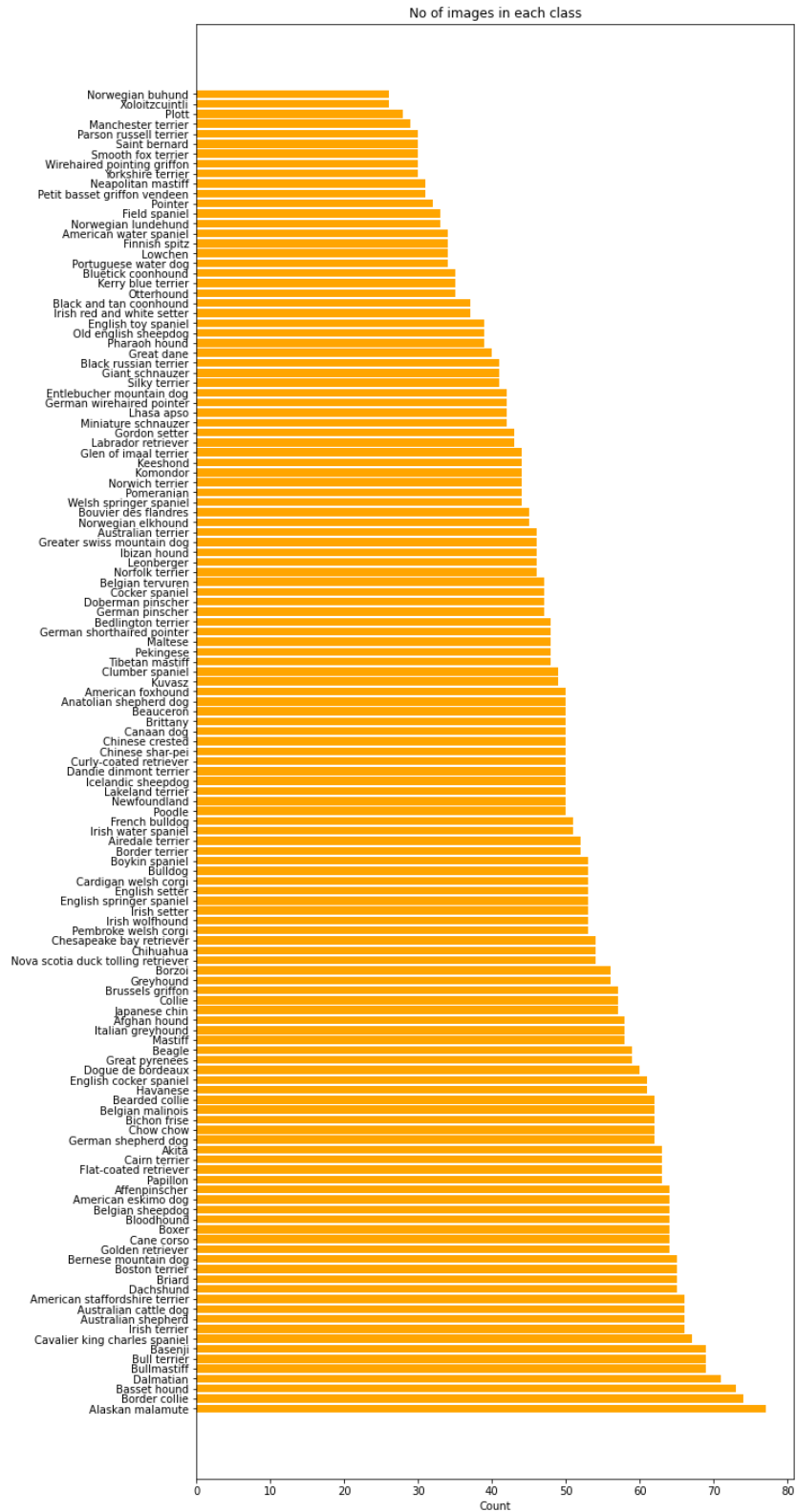


Fig. 2 Bar chart showing total count of images from each class

Algorithms and Techniques

Our goal in this project is to create a CNN from scratch as well as by using transfer learning to classify dog breeds. To start, the CNN receives an input feature map: a three-dimensional matrix where the size of the first two dimensions corresponds to the length and width of the images in pixels. The size of the third dimension is 3 (corresponding to the 3 channels of a color image: red, green, and blue). The CNN comprises a stack of modules, each of which performs three operations: Convolution [6], Rectified Linear Unit [7] (ReLU) and Pooling [8].

At the end of a convolutional neural network are one or more fully connected layers [9] (when two layers are "fully connected," every node in the first layer is connected to every node in the second layer). Their job is to perform classification based on the features extracted by the convolutions. Here we use a final fully connected layer having 133 output features to classify the input image into one of the 133 dog breed classes.

As with any machine learning model, a key concern when training a convolutional neural network is overfitting [10]: a model so tuned to the specifics of the training data that it is unable to generalize to new examples. Two techniques to prevent overfitting when building a CNN are: Data augmentation and Dropout regularization. Here we use Dropout to prevent overfitting.

Training a convolutional neural network to perform image classification tasks typically requires an extremely large amount of training data, and can be very time-consuming, taking days or even weeks to complete. But what if you could leverage existing image models trained on enormous datasets, such as Inception, and adapt them for use in your own classification tasks? One common technique for leveraging pretrained models is feature extraction: retrieving intermediate representations produced by the pretrained model, and then feeding these representations into a new model as input.

Here we use a pretrained ResNet50 [12] model to create a CNN based on transfer learning. Here we replace the final fully connected layer of the pretrained ResNet50 model with our own having 133 output features corresponding to the 133 dog breeds.

Benchmark

The CNN model built from scratch by us will serve as a benchmark model to validate the performance of our final optimized model built using transfer learning. This benchmark model also ascertain that the dog breed classification problem is clearly defined and solvable.

Methodology

Data Preprocessing

The dog dataset was already separated into train, valid and test datasets and each dataset was further separated into 133 folders corresponding to the 133 dog breeds.

We applied the following preprocessing and data augmentation techniques:

- Randomly resized and cropped the images to 224x224 using *transforms.RandomResizedCrop(224)*
- Randomly flipped the images horizontally using *transforms.RandomHorizontalFlip()*
- Randomly rotated the images using *transforms.RandomRotation(10)*
- Normalized the images using mean = [0.485, 0.456, 0.406] and std = [0.229, 0.224, 0.225]

Implementation

The implementation consists of the following steps:

1. **Import Datasets:** Here we download and import the required human and dog datasets.
2. **Detect Humans:** Here we use OpenCV's implementation of Haar feature-based cascade classifiers to detect human faces in images.

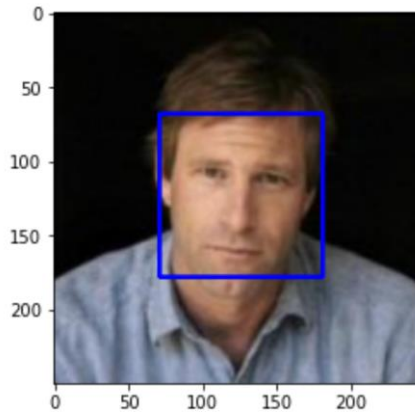


Fig.3 Face detection using [Haar feature-based cascade classifiers](#)

We have tested the Human Detector algorithm on both the human and dog datasets.

The results are as follows:

```
Human Images with detected faces - 96.0%  
Dog Images with detected faces - 18.0%
```

3. **Detect Dogs:** Here we use a pre-trained VGG-16 model to detect dogs in images. We have downloaded the VGG-16 model, along with weights that have been trained on ImageNet, a very large, very popular dataset used for image classification and other vision tasks. ImageNet contains over 10 million URLs, each linking to an image containing an object from one of 1000 categories. We need only check if the pre-trained model predicts an index between 151 and 268 (inclusive), to include all categories from *Chihuahua* to *Mexican hairless*.

We have tested the Dog Detector algorithm on both the human and dog datasets.

The results are as follows:

```
Human Images with detected dogs - 0.0%  
Dog Images with detected dogs - 95.0%
```

The Dog Breed Classification consists of the following two stages:

1. **Create a CNN to Classify Dog Breeds (from Scratch):**

Here we will create a CNN that classifies dog breeds from scratch and use it as a benchmark model to evaluate the CNN model which we create using Transfer Learning.

We have chosen the following architecture for our CNN model:

```
Net(  
    (conv1): Conv2d(3, 6, kernel_size=(3, 3), stride=(1, 1))  
    (conv2): Conv2d(6, 12, kernel_size=(3, 3), stride=(1, 1))  
    (conv3): Conv2d(12, 24, kernel_size=(3, 3), stride=(1, 1))  
    (conv4): Conv2d(24, 48, kernel_size=(3, 3), stride=(1, 1))  
    (fc1): Linear(in_features=6912, out_features=1024, bias=True)  
    (fc2): Linear(in_features=1024, out_features=512, bias=True)  
    (fc3): Linear(in_features=512, out_features=133, bias=True)  
)
```

The CNN architecture consists of 4 Convolutional layers and 3 Fully Connected layers. ReLU activation function was applied to every layer except the last one. MaxPool2d was applied to the Convolutional layers to reduce the output into half. An input image of size 224x224x3 is passed and output of 133 is obtained.

1st Layer: Consists of Conv2d with 3 Input Channels and 6 Output Channels and a Kernel size of 3x3 and stride 1x1. It is followed by ReLU and Maxpool2d.

Here an input of size 224x224x3 is passed and output of 111x111x6 is obtained.

2nd Layer: Consists of Conv2d with 6 Input Channels and 12 Output Channels and a Kernel size of 3x3 and stride 1x1. It is followed by ReLU and Maxpool2d.

Here an input of size 111x111x6 is passed and output of 54x54x12 is obtained.

3rd Layer: Consists of Conv2d with 12 Input Channels and 24 Output Channels and a Kernel size of 3x3 and stride 1x1. It is followed by ReLU and Maxpool2d.

Here an input of size 54x54x12 is passed and output of 26x26x24 is obtained.

4th Layer: Consists of Conv2d with 24 Input Channels and 48 Output Channels and a Kernel size of 3x3 and stride 1x1. It is followed by ReLU and Maxpool2d.

Here an input of size 26x26x24 is passed and output of 12x12x48 is obtained.

5th Layer: Consists of Linear with 6912 Input Features and 1024 Output Features.

6th Layer: Consists of Linear with 1024 Input Features and 512 Output Features.

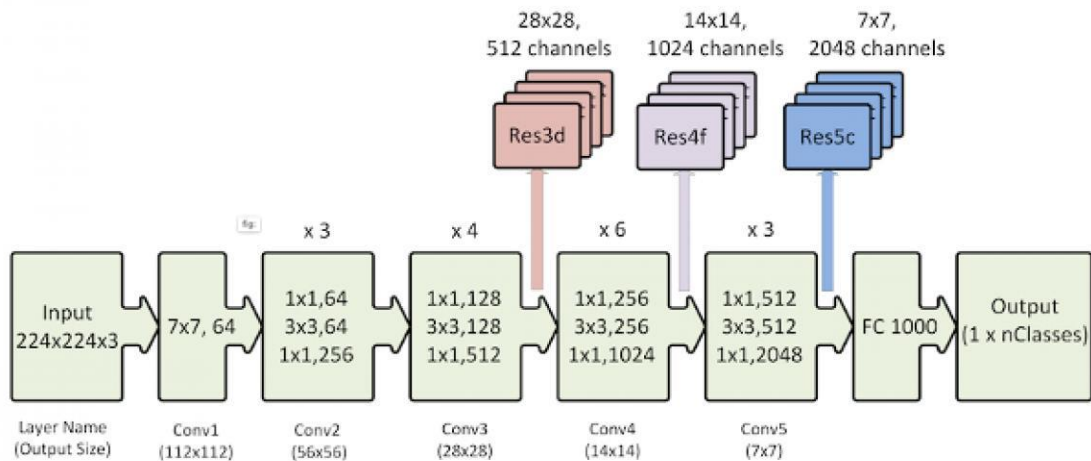
7th Layer: Consists of Linear with 512 Input Features and 133 Output Features.

We have chosen Cross Entropy Loss as our loss function and Adam optimizer and a learning rate of 0.001. We have trained the model for 20 epochs and achieved a test accuracy of 16%.

2. **Create a CNN to Classify Dog Breeds (using Transfer Learning):**

We will now use transfer learning to create a CNN that can identify dog breed from images. We used a pretrained ResNet50 model from Torchvision Model Zoo for our Dog Classification problem. This pretrained model was trained on ImageNet dataset and has the last Fully Connected layer with 1000 out features. These pre-trained networks demonstrate a strong ability to generalize to images outside the ImageNet dataset via transfer learning.

Architecture diagram:



Refinement

We made modifications in the pre-existing model by fine-tuning the model. We have frozen the parameters, so we don't back-propagate through them and replaced the last Fully Connected layer with a Linear layer having 133 out features equal to the number of classes in our Dog dataset. Since we assume that the pre-trained network has been trained quite well, I think this architecture is suitable for our current problem after fine-tuning the model.

We have chosen Cross Entropy Loss as our loss function and Adam optimizer and a learning rate of 0.001. We have trained the model for 20 epochs and achieved a test accuracy of 85%.

Here is the output of the training cycle:

```
Epoch: 1      Training Loss: 2.740375      Validation Loss: 0.884047
Validation loss decreased from inf to 0.884047. Saving model ...
Epoch: 2      Training Loss: 1.381115      Validation Loss: 0.658374
Validation loss decreased from 0.884047 to 0.658374. Saving model ...
Epoch: 3      Training Loss: 1.170744      Validation Loss: 0.584149
Validation loss decreased from 0.658374 to 0.584149. Saving model ...
Epoch: 4      Training Loss: 1.054523      Validation Loss: 0.540966
Validation loss decreased from 0.584149 to 0.540966. Saving model ...
```

Epoch: 5	Training Loss: 1.009576	Validation Loss: 0.514517
Validation loss decreased from 0.540966 to 0.514517. Saving model ...		
Epoch: 6	Training Loss: 0.972923	Validation Loss: 0.479021
Validation loss decreased from 0.514517 to 0.479021. Saving model ...		
Epoch: 7	Training Loss: 0.923414	Validation Loss: 0.492444
Epoch: 8	Training Loss: 0.920773	Validation Loss: 0.510267
Epoch: 9	Training Loss: 0.925606	Validation Loss: 0.501688
Epoch: 10	Training Loss: 0.884902	Validation Loss: 0.465351
Validation loss decreased from 0.479021 to 0.465351. Saving model ...		
Epoch: 11	Training Loss: 0.840818	Validation Loss: 0.482059
Epoch: 12	Training Loss: 0.901755	Validation Loss: 0.432933
Validation loss decreased from 0.465351 to 0.432933. Saving model ...		
Epoch: 13	Training Loss: 0.874118	Validation Loss: 0.479497
Epoch: 14	Training Loss: 0.844539	Validation Loss: 0.497694
Epoch: 15	Training Loss: 0.863298	Validation Loss: 0.529868
Epoch: 16	Training Loss: 0.842218	Validation Loss: 0.502431
Epoch: 17	Training Loss: 0.827928	Validation Loss: 0.496574
Epoch: 18	Training Loss: 0.805357	Validation Loss: 0.463499
Epoch: 19	Training Loss: 0.815395	Validation Loss: 0.463986
Epoch: 20	Training Loss: 0.826752	Validation Loss: 0.497854

Clearly, we saw that our new model we created using transfer learning outperformed the previous model created from scratch. The training/validation loss is decreased while the test accuracy is increased.

Results

Model Evaluation and Validation

Here we evaluated our Resnet50 model against the test dataset and calculated the test accuracy using the following formula:

$$\text{Accuracy} = (\text{True Positives} + \text{True Negatives}) / \text{Total Images}$$

We achieved a test accuracy of 85%.

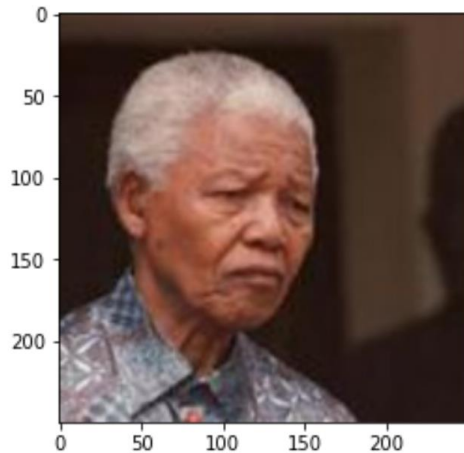
Next, we created a Dog Breed detection pipeline using the above model that accepts a file path to an image and first determines whether the image contains a human, dog, or neither.

Then,

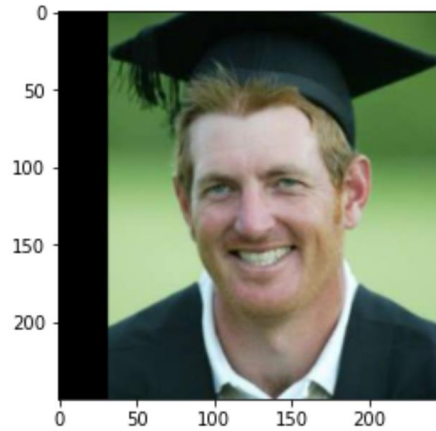
- if a dog is detected in the image, return the predicted breed.
- if a human is detected in the image, return the resembling dog breed.
- if neither is detected in the image, provide output that indicates an error.

We then tested our pipeline on few images from the both datasets and the results were satisfactory. Here are some of the detections of our Dog Breed Classification pipeline.

Hello Human!
You look like a...
American water spaniel



Hello Human!
You look like a...
Doberman pinscher



Hello Dog!
Your predicted breed is...
Bull terrier



Hello Dog!
Your predicted breed is...
Italian greyhound



Justification

We compared the test accuracy scores of both the CNN model created from scratch and the CNN model created using Transfer Learning and found that the final model which used Transfer Learning outperformed the previous model created from scratch. The new model achieved an accuracy score of 85% while the old model achieved merely 16% when run for the same number of epochs.

While the final model is significant enough to have adequately solved the problem, there is always some room for improvements. We can try the following techniques to improve the test accuracy for our final model:

- Perform histogram equalization to augment the dataset
- Try a different Optimizer like SGD and lower the learning rate.
- Use dropouts to prevent overfitting
- Use more layers like Resnet101 or Resnet152

References

1. <https://www.sciencedirect.com/topics/computer-science/image-classification>
2. https://en.wikipedia.org/wiki/Supervised_learning
3. https://wikipedia.org/wiki/Convolutional_neural_network
4. <https://s3-us-west-1.amazonaws.com/udacity-aind/dog-project/lfw.zip>
5. <https://s3-us-west-1.amazonaws.com/udacity-aind/dog-project/dogImages.zip>
6. https://en.wikipedia.org/wiki/Convolutional_neural_network#Convolutional_layer
7. https://en.wikipedia.org/wiki/Convolutional_neural_network#ReLU_layer
8. https://en.wikipedia.org/wiki/Convolutional_neural_network#Pooling_layer
9. https://en.wikipedia.org/wiki/Convolutional_neural_network#Fully_connected_layer
10. <https://en.wikipedia.org/wiki/Overfitting>
11. <https://arxiv.org/abs/1409.1556>
12. <https://arxiv.org/abs/1512.03385>