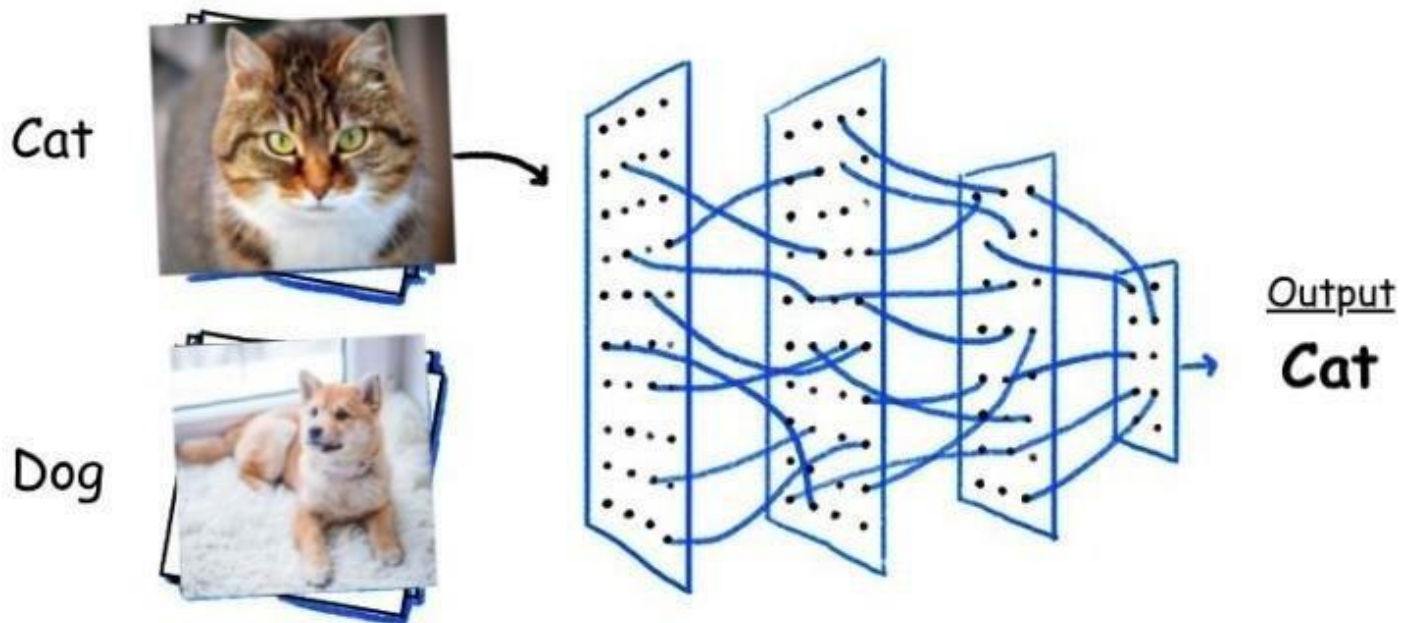


Cats Vs Dogs Identification Using CNN



What Is A CNN Model?

In deep learning, a convolutional neural network (CNN/ConvNet) is a class of deep neural networks, most commonly applied to analyze visual imagery. Unlike the classical image recognition where you define the image features yourself, CNN takes the image's raw pixel data, trains the model, then extracts the features automatically for better classification.

There are various types of CNN architecture like AlexNet, VGGNet, ResNet, and GoogleNet.

CNN is mainly composed of an Input layer, Convolution layers, Activation function (ReLU), Pooling layer, Flattening, Fully connected layer, and Softmax or Sigmoid function depending on whether you're doing multiclass or binary classification.

The Input Layer

The input layer is the first layer in any CNN model. It is the layer that receives the image for the CNN model. The dimension of the image along with the image channels is passed to the input layer.

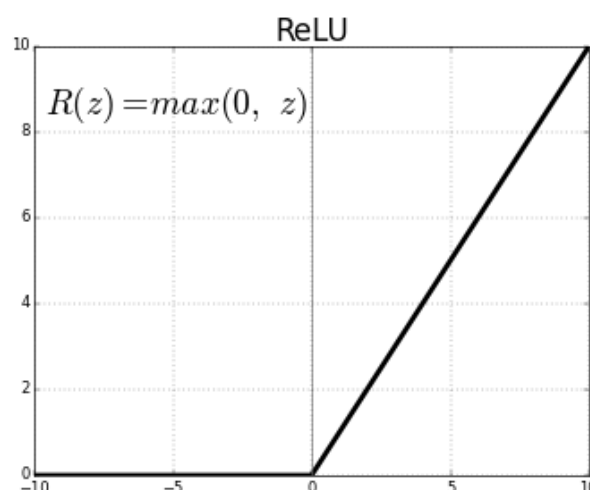
The Convolution Layer

This is the first layer to extract features from the image. These extractions are done with the help of filters. It does a mathematical operation on the image with the filter and preserves the relationship between pixels by learning image features.

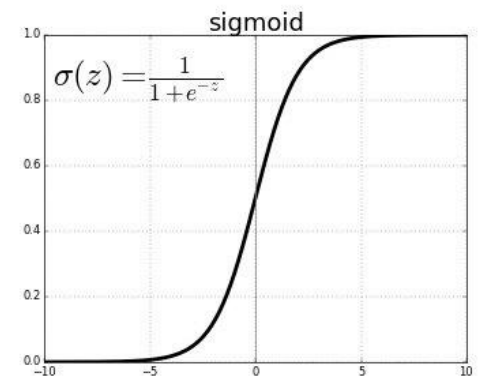
Activation Function

The activation function is a node that is put at the end of or in between Neural Networks. They help to decide if the neuron would fire or not. The most popular activation functions in the CNN model are ReLU, Sigmoid, and Softmax.

ReLU:- The rectified linear activation function or ReLU for short is a piecewise linear function that will output the input directly if it is positive, otherwise, it will output zero.



Sigmoid:- This function is mostly used for the model which has binary outputs and is added to the output layer. It exists between 0 and 1.



Softmax:- This Function is also added in the output layer where we have multiple classes to predict. Softmax assigns decimal probabilities to each class in a multi-class problem. Those decimal probabilities must add up to 1.0

The Pooling Layer

Pooling layers are used to reduce the dimensions of the feature maps. It reduces the number of trainable parameters and hence the number of computations required. There are mainly two types of pooling layers used MaxPool and Average Pool.

Flattening

Flattening is converting the data into a 1-dimensional array for inputting it to the next layer. We flatten the output of the convolutional layers to create a single long feature vector that can be passed to fully connected layers.

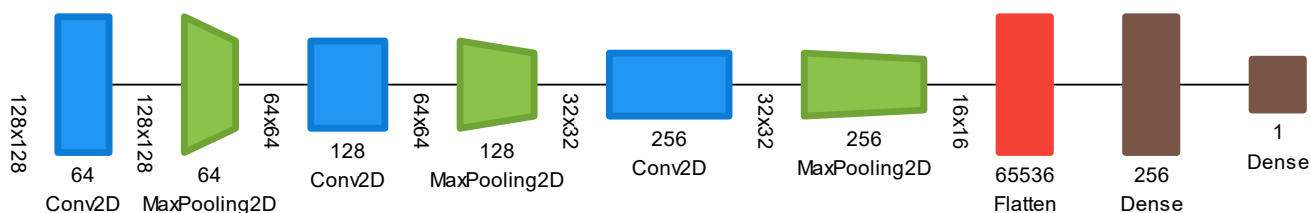
Fully Connected Layers

Fully Connected Layers form the last few layers in the network. The input to the fully connected layer is the output from the final Pooling or Convolutional Layer, which is flattened and then fed into the fully connected layer.

Cat-vs-Dogs Classification Model Architecture

The Custom Made Cat-vs-Dog Classification Model consists of 3 convolutional layers with 3 max-pooling layers and 1 fully connected layer.

Its architecture is given below.



Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 128, 128, 64)	4864
max_pooling2d (MaxPooling2D)	(None, 64, 64, 64)	0
conv2d_1 (Conv2D)	(None, 64, 64, 128)	73856
max_pooling2d_1 (MaxPooling2D)	(None, 32, 32, 128)	0
conv2d_2 (Conv2D)	(None, 32, 32, 256)	295168
max_pooling2d_2 (MaxPooling2D)	(None, 16, 16, 256)	0
flatten (Flatten)	(None, 65536)	0
dense (Dense)	(None, 256)	16777472
dense_1 (Dense)	(None, 1)	257
Total params: 17,151,617		
Trainable params: 17,151,617		
Non-trainable params: 0		

Deep Learning Concepts Used

Neural Network

It is formed by interconnected neurons. These neurons have weights, and bias which is updated during the network training. The network consists of an input layer, hidden layers, and an output layer. The hidden layer is where the processing of data is done.

Activation Functions

The activation function translates the input signals to output signals. We apply a nonlinear function like an activation function to a linear function coming from the input layer. Common activation functions are sigmoid, Relu, Softmax, and tanh.

Loss Functions

When we build a network, the network tries to predict the output as close as possible to the actual value. We measure this accuracy of the network using the loss function. Our motive while running the network is to increase our prediction accuracy and to reduce the error, hence minimizing the cost function.

Gradient Descent

Gradient descent is an optimization algorithm for minimizing the loss. The basic idea behind it is that the gradient of the function should try to find the direction of the steeper section (towards minima).

Batches

While training a neural network, instead of sending the entire input in one go, we divide the input into several parts of equal size randomly. Larger batch sizes help in reducing the training time along with the cost of little low accuracy.

Epochs

An epoch is defined as a single training iteration of all batches in both forward and backpropagation. This means 1 epoch is a single forward and backward pass of the entire input data.

Batch Normalization

Batch Normalization is done to ensure that the distribution of data is the same as the next layer hoped to get. When we are training the neural network, the weights are changed after each step of gradient descent. It is performed on the output of layers of each batch.

Dataset

The dataset consists of 4000 images of cats and dogs each for training and 1000 images of both for test in directories. The image sizes of all images are different and are colored images.

These datasets are loaded from the directory into the Training and Testing set and a part of the training set (11%) will be used as a validation set while training. Further, it is divided into the batch size of 32, and the image size is rescaled to 128 by 128 with the help of data augmentation.

Data Augmentation

It is a commonly used technique in image processing used to create more training data. It modifies the original training images a bit for example by rotating, flipping, zooming in and out, etc.

Image Pre-Processing

Image preprocessing is the steps taken to format images before they are used by model training and inference. This includes resizing, orienting, and color corrections.

A few of them which are used in my image is explained below.

Rescaling

Every digital image is formed by pixels having values in the range 0~255. 0. Similarly, in color images also for all three channels, the range is from 0~255.

Since 255 is the maximum pixel value. Rescale $1./255$ is to transform every pixel value from range $[0,255] \rightarrow [0,1]$.

This is done because some images are in a high pixel range, and some are in a low pixel range. The images are all sharing the same model, weights, and learning rate. So if we don't scale it equally some differences will occur in the images and either high loss or weak loss is generated. Also as training neural networks with this range gets efficient

Resizing

Resizing is used as the images in the data set are of different dimensions and some of them are quite large. The larger the image the more data can be gathered by training time will also increase tremendously. That's why resizing each image to a small size is necessary. Images are generally resized to 32x32, 128x128, 256x256 or 512x512.

In my model, I have resized it to 128x128.

Project Details

In this project, I have built a CNN model from scratch and trained it with the help of the dataset from [here](#). First I pre-processed the image by rescaling it and applying some image augmentations like horizontal flip, rotation by 40 deg, and also resized it to 128x128. And then I split the training dataset into training and validation sets. All this was done with the help of ImageDataGenerator from the keras.preprocessing library.

Then I created a simple CNN model with 3 Convolution layers each followed by a max-pooling layer. Each layer had ReLU as its activation function except the output layer which had a sigmoid.

Then I compiled the model with adam as optimizer and binary_crossentropy as the loss function.

For training, I also used callbacks and early stopping as I didn't know how much longer it will take to train my model. With a batch size of 64, the model was trained for 25 epochs and the training accuracy came around 75%, final validation accuracy around 77%, and test accuracy around 80%. The training took around 25 minutes to complete with the help of Kaggle's GPU.

Conclusion

In this project, an ML model was built to identify an image as a cat or a dog with the help of a convolutional neural network. The model's final accuracy was around 80% in around 25 minutes of training with twenty-five epochs.