

Deep learning - Dog Breed Classification

Realization of an native Android app using deep learning algorithms

Alice Bollenmiller, Andreas Wilhelm
WS 17/18 IG
January 26, 2018

Contents

1. Introduction	1
1.1. Deep learning	1
1.2. Terms of Referencee	1
2. Methodological fundamentals	1
2.1. Common Frameworks for Deep Learning Applications	2
2.1.1. Tensorflow	2
2.1.2. Keras	2
2.1.3. Caffé	2
2.1.4. Torch and PyTorch	2
2.2. Common Models in Deep Learning Applications	2
2.3. Key requirements for an appropriate dataset	6
3. Concept	7
3.1. Frameworks	7
3.2. Model based Architectures	7
3.3. Application based Architecture	7
4. Realization	9
4.1. dataset	9
4.2. hardware environment	9
4.3. installation of software	9
4.3.1. Prerequisites	9
4.3.2. Tensorflow based on Python	10
4.3.3. Tensorflow based on Bazel	10
4.3.4. Installing Android Studio and its Delevopment Kit	10
4.4. building the models	10
4.5. Output Tests and Validation	10
4.6. Implementation of an native Android App	10
4.7. Deployment and Validation	11
5. Evaluation	12
6. Conclusion	12
Bibliography	I
List of figures	I
A. Anhang	II
A.1. Anhang vom Anhang	II

1. Introduction

In this chapter, a short introduction leads to the subject of Deep learning. Furthermore, the scope of this work and its points of reference are described and localized.

1.1. Deep learning

In 2015, AlphaGo - a computer program developed by Google's DeepMind Group that was trained to play the strategic board game Go - was the first program to defeat the multiple European champion Fan Hui under tournament conditions. Back then, terms like Artificial Intelligence (AI), Machine Learning and Deep Learning were talked of, because those were the reason why a computer program was able to defeat a human being. One of the most frequent used techniques of AI is Machine Learning. It uses algorithms to parse data, process it and learn from it. The result is a prediction or determination as a conclusion of what was learnt from the dataset. Deep Learning - which is part of the Machine Learning techniques - sells its application particularly in the field of language and image processing.

In 1958, Frank Rosenblatt introduced the concept of the perceptron which is the fundamental idea of all Deep learning approaches. It consists of multiple artificial neurons which are coupled with weights and biases. In the case of a single-layer perceptron, the input nodes are fully connected to one or more output nodes. During the learning process the weights are adapted according to the learning progress. When such a structure is extended with layers, it becomes a multi-layer perceptron (MLP). This basic structure can be found in special neural network architectures e.g. Convolutional Neural Networks (CNN). CNNs which are frequently used for object detection and image/ audio recognition etc. consists of multiple neurons. When a neuron receives an input, it performs a dot production and adapts the weights. Because of their special structure, CNNs are able to detect local properties of an image. Basically, the network represents a differentiable score function which is applied on the raw image pixels and computes the class scores as an output.

1.2. Terms of Reference

The problem of Deep learning architectures is their high performance requirements regarding computational power. Common frameworks for open source development in the field of deep learning which are discussed in section 2.1 introduced several models for the integration in mobile applications.

In order to overcome the above mentioned problem, optimizing methods will be combined with appropriate models which require less computing power and are suitable for mobile application integration. As a result of this work, a dog breed analyzer will be implemented. This mobile app will take a live camera stream as an input and determine the breed of the focused dog. The three highest probabilities of breeds will be shown by the app.

2. Methodological fundamentals

This chapter describes the most frequently used frameworks in deep learning for developing applications. Furthermore, common models for deep learning are introduced followed by suitable models for mobile integration. The chapter closes listing key requirements for an appropriate dataset which increase the quality of the training results.

2.1. Common Frameworks for Deep Learning Applications

The demands on neural networks increases with the complexity of problems to solve. Concurrently, there's an expanding offer of deep learning frameworks with a variety of features and tools. The most common used ones are represented in the following section.

2.1.1. Tensorflow

In 2015, the Google Brain Team introduced the most popular deep learning API Tensorflow which is an open-source library for numerical computation. Its current version 1.4.1 was released on December 8th, 2017. Tensorflow is primarily used for machine learning and deep neural network research. Based on the programming language Python, Tensorflow is capable of running on multiple CPUs and GPUs. Furthermore, C++ and R are supported by Tensorflow. Another feature is the possibility to generate models and export them as .pb file which holds the graph definition (GraphDef). The export is done by protocol buffers (protobuf) which includes tools for serializing and processing structured data. When loading a .pb file by protobuf, a graph object is created which holds a network of nodes. Each of those nodes represent an operation and the output is used as input for another operation. This concept enables an user to create self-built tensors. To simplify the usability, Tensorflow developed a high-level wrapper of the native API which is called Tensorflow Slim. Furthermore, in order to run Tensorflow on performance critical devices like e.g. mobile devices there are two lightweight solutions of Tensorflow available: Tensorflow Mobile and Tensorflow Lite. The latter one is an evolution of Tensorflow Mobile and still in developer mode. But both are predestinated for integration in mobile applications.

2.1.2. Keras

In order to simplify the utilization of Tensorflow the Python based interface Keras can be configured to work on top of Tensorflow. It allows building neural networks in a simple way and is part of Tensorflow.

2.1.3. Caffe

Another deep learning library is Caffe which was developed by Berkeley AI Research (BAIR). Based on C++ or Python, it focuses on modeling CNNs. A main advantage of Caffe is the offer of pretrained models available in the Caffe Model Zoo.

2.1.4. Torch and PyTorch

Besides Tensorflow and Caffe, Torch is another common deep learning framework. It was developed by Facebook, Twitter and Google. Based on C/C++, Torch supports CUDA for GPU processing. Like above mentioned frameworks, Torch facilitates the building of neural networks. The Python based version of Torch is available through PyTorch.

2.2. Common Models in Deep Learning Applications

State of the art models for image recognition in deep learning applications are in a CNNs architecture. The most important ones are explained in this section.

The first one was LeNet-5 of ?. As exactly shown in Figure 1 it comprises seven layers. Thereof two times a convolutional layer with a 5x5 filter each followed by a sub-sampling pooling layer. Two fully connected layers complete the architecture of the multilayer perceptron with its softmax function applied.

In 2012 a CNN based network called AlexNet of ? won the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) for the first time (Figure 2). From that time CNN became ubiquitous.



Figure 1: Model architecture of LeNet-5 (?)



Figure 2: ImageNet Classification top-5 error (%) (?)

Through its immensely improved accuracy the network enabled a lot new possibilities. This CNN now consists of eight layers shown in Figure 3. The first convolutional layer has a large 11x11 filter. Furthermore one with 5x5 and another three with a 3x3 filter, so in total 5 convolutional layers. Three pooling layers and also three fully connected layers completes the architecture. For the first time a rectified linear unit (ReLU) was used and it could be spread across two GPUs.



Figure 3: Model architecture of AlexNet (?)

As seen in Figure 2 the next groundbreaking improvements were two architectures in 2014 which were also considered as very deep networks. VGG of ? and GoogLeNet of ? which is today better known under 'Inception'.

The first one comprises 16 or 19 layers in total (Figure 4). This network is characterized by its simplicity because it uses convolutional layers with a 3x3 filter which only stacked one above the other. Reducing the size is handled by pooling layers. At the end there are also three fully connected layers with the softmax function following.



Figure 4: Model architecture of VGG19 (?)

The other one has 22 layers (Figure 5) and was the final winner in 2014. GoogLeNet starts with two convolutional layers followed by stacked inception modules. Such a module consists of several parallel convolutional layers with different filter size and also one pooling layer as shown in Figure 6. On the right side of this figure a convolutional layer with a 1x1 filter is used to reduce feature depth. On the left side all output of these modules are concatenate and forwarded as new input for the next layer. The network ends with only one fully connected layer.

Since 2015 it gets very deep with a new architecture named ResNet of ?. On ILSVRC this network swept all competitors in classification and detection. The network itself is designed as a 152 layer model. The ResNet consists of many residual network blocks (Figure 7). Each of them has two convolutional layers with a 3x3 filter. Furthermore there is a shortcut connection which will be



Figure 5: Model architecture of Inception (?)



Figure 6: Inception modules (?)

used if the input and output dimension of this block are the same.



Figure 7: ResNet block (?)

Especially for performance critical applications like mobile devices MobileNet was developed by Google respectively ?. The network structure consists of depthwise seperable convolutional parts but also one standard convolutional one at the beginning. The depthwise part is separated in a depthwise convolutional layer and a pointwise one which applies a 1x1 convolution. Each of this layers are followed by batchnormalisation and ReLu as shwon in Figure 8. In the depthwise convolution the splitting of the two convolutional layers in one for filtering and the other one for combining has the effect to reduce massively computation and model size which results in better performance on low performance devices.



Figure 8: Standard convolution and depthwise seperable convolution (?)

2.3. Key requirements for an appropriate dataset

Supervised learning tasks such as image recognition are based on operations where an output is taken as an input for the next node. Every raw pixel input is taken to compute an intermediate representation - a vector containing all learned information about the dataset. As a consequence, the training results are only as good as the dataset itself. For better accuracy its important to train a model on a variety of images for each object which should later be classified by the model. It's recommended to take images of an object which were taken at different times, with different devices and at different places. Otherwise, the model will concentrate on other things like for example the background instead of details about the object itself. Therefore, a huge dataset is required especially for non pre-trained models. Training a model from scratch will require a huge dataset, a lot computing power and time. Whereas pre-trained models only require a small dataset

of about hundreds of images. For that reason, a pre-trained model will be used in this work.

3. Concept

First, this chapter describes the selection of the appropriate framework. Furthermore, the structure of the model which was used for classification is explained based on its architecture. The chapter closes with the class diagram of the mobile application.

3.1. Frameworks

As a deep learning framework Tensorflow was used to retrain the model. This decision was mainly based on recommendations. Even companies like e.g. NVIDIA Corporation, Intel Corporation etc. use Tensorflow. It is one of the common frameworks for deep learning applications and also provides solutions for integration in mobile apps. Furthermore, Tensorflow provides a variety of tutorials for working with neural networks. Beside those advantages, there is a large community about Tensorflow talking about issues and solutions.

To run the tensor within a mobile application, the first approach was to use Tensorflow Lite which is still in development state. But many attempts resulted in corrupt models which caused the app to terminate. Because of this experiences Tensorflow Mobile was used to run the model within the mobile application.

3.2. Model based Architectures

Andi - general architectures of models - n Mobilenet, Inception

3.3. Application based Architecture

This chapter describes the functions to a basic understanding of how the application works based on the application's architecture. It focuses on the important classes and methods of the mobile application.

Tensorflow provides a mobile application for demo purposes. Because of the complexity the application was adapted to our needs. The approach was to understand how the application was implemented regarding the functionality of the camera, image input for the network and background classification besides live camera stream.

The ClassifierActivity is set as the launcher activity of the application. It extends the class CameraActivity and loads the CameraFragment which controls the camera view. The CustomTextureView enables the possibility to capture frames from a camera stream and process it. Because of this function, the CustomTextureView is part of the class CameraFragment. If an image is captured, the method onImageAvailable of the class ClassifierActivity is invoked. The classification itself is done by the TensorflowImageClassifier which implements the interface Classifier from the Tensorflow API. The results are displayed by the class RecognitionScoreView which implements the interface ResultsView.

4. Realization

In general, this chapter describes the methodical procedure of solving the above mentioned problem section 1.2. After describing the used dataset all required software and hardware components are explained in detail. Furthermore, the chapter leads through the installation steps of Tensorflow

and the setup of Android Studio. Followed by the installation process the retraining of a pre-trained model is depicted. Afterwards, the re-trained model is tested and validated. The chapter ends with the description of the realization of the Android app.

4.1. dataset

Andi

4.2. hardware environment

Andi used CPU, GPU -i NVIDIA, handys

4.3. installation of software

Andi This chapter includes all necessary steps for installing the software environment including Tensorflow.

4.3.1. Prerequisites

The software environment was set up on the Linux distribution Ubuntu 16.04 LTS. To install the software environment for Tensorflow Python is required. Therefore, the current version of Python 3.6 was installed by default. Tensorflow also supports Bazel which was installed by following command.

```
1 sudo apt-get install openjdk-8-jdk
2
3 echo "deb [arch=amd64] http://storage.googleapis.com/bazel-apt stable jdk1.8" | sudo
   tee /etc/apt/sources.list.d/bazel.list
4 curl https://bazel.build/bazel-release.pub.gpg | sudo apt-key add -
5
6 sudo apt-get update && sudo apt-get install bazel
7
8 sudo apt-get upgrade bazel
```

Listing 1: Bazel Installation

Futhermore, the package and environment management tool Anaconda was installed by the following steps: First, the Anaconda installer was downloaded from <https://www.anaconda.com/download/#linux>. During the installation process the prompts were answered by the default suggestions except the following prompt: "Do you wish the installer to prepend the Anaconda3 install location to PATH in your /home/aw/.bashrc ? [yes—no]". "yes" was typed in and conda was tested using the "conda list" command.

As an app development environment the free IDE Android Studio in its version 3.0.1 was installed by downloading from <https://developer.android.com/studio/index.html>, extracting and following the instruction steps. Required dependencies are installed by Android Studio itself. So, the SDK in the version 26.1.1 was used.

In the development of native Android Apps Java is used as the programming language. Within the installation of Anaconda, the JDK in the version 8 was installed. - test anaconda, tensorflow environment look where it belongs

- CUDA, CUDNN

4.3.2. Tensorflow based on Python

4.3.3. Tensorflow based on Bazel

- e.g. Workspace changes for Android SDK, msse4.2

4.3.4. Installing Android Studio and its Development Kit

- also possible with bazel but easier Android studio (needs correct versions of sdk, ndk)
- SDK, NDK
- IMPORTANT: tf versions updaten (same as trained)

4.4. building the models

Alice bis Steps, Andi ab Optimierung, time GPU/CPU -> evtl extra subsection:

- execution methods -> Bazel and Python (incompatible versions)
- Mobilnet -> steps, optimierung
- Inception -> steps, optimierung
- time related differences of execution
- > time CPUs/GPU

4.5. Output Tests and Validation

- Andi - test pictures and if it works -> label image
- validation script?! -> in Evaluierung

4.6. Implementation of an native Android App

Alice This subchapter describes the processing of the camera input stream and classifying the particular images. Because of the vast extent of the application, the focus lies on the explanation of how the app works and on its important implementation parts.

When click on the icon of the app, a camera view showing results in the bottom part is loaded. In the background, the ClassifierActivity is set as the launcher activity in the AndroidManifest.xml which is a file containing all configuration settings for the app. When an activity is loaded for the first time, the onCreate method is invoked. Because the ClassifierActivity extends the CameraActivity the onCreate method of the CameraActivity is executed and loads the layout activity_camera.xml consisting of a container and a result field. Next, the setFragment method of the class CameraActivity is invoked if permission for the camera is given. This causes the replacement of the container with the adapted fragment according to the actual size and orientation of the screen. When the size is set, the method onPreviewSizeChosen is invoked where a TensorFlow-ImageClassifier is instantiated with following input parameters: MODEL_FILE, LABEL_FILE, INPUT_SIZE, IMAGE_MEAN, IMAGE_STD, INPUT_NAME, OUTPUT_NAME. Details about the Classification follow later in the section.

When the fragment was replaced by the CameraActivity, an instance of the class CameraFragment was created. Next, the method onViewCreated of the class CameraFragment is called and a texture view is created. The latter class is part of the Android API and used for frame capturing from an image stream as OpenGL ES texture. Thereby, the most recent image is stored within a texture object. This object is observed by a listener which invokes the method onSurfaceTextureAvailable when the texture object is available. Within this method the camera is opened by the method openCamera(width, height) given the width and height of the camera preview (?). First, the method openCamera sets the camera parameters within setUpCameraOutputs e.g. sensorOrientation, previewSize, cameraId etc. In order to classify a given OpenGL ES texture, the coordinate column vectors of the texture must be transformed into the proper sampling location in the streamed texture. So, the matrix has to be prepared with the correct configuration which happens in the method configureTransform. The next important step is done by the camera manager which opens the camera by invoking the method openCamer(cameraId, stateCallback, backgroundHandler). The id represents the specific camera device whereas the stateCallback is necessary to manage the life

circle of the camera device and to handle different states of the camera device. The background-Handler ensures that the classification is done in background mode. When the camera is opened, the camera preview is started by the method `createCameraPreviewSession`. Within this method the preview reader is initialized and set to read images from the `ImageListener` which observes whether an image is available. Therefore, an `ImageListener` which was instantiated by the `CameraActivity` is transferred to the `CameraFragment`.

When an image is available the method `onImageAvailable` of the class `ClassifierActivity` is invoked. The image is read by the preview reader and stored in an object. First, the image is preprocessed meaning transferred into planes, stored as bytes, cropped and stored as a `Bitmap` which is an Android graphic. Google provides an `Tensorflow Mobile API` to run tensors on a performance critical device such as mobile devices. To use this API, the dependency has to be added to the `build.gradle` file listing 4.

```

1
2 dependencies {
3     compile 'org.tensorflow:tensorflow-android:+'
4 }

```

Listing 2: Tensorflow API in build.gradle

When including the API, Tensorflow's class `Classifier` can be used to classify images. Therefore, the `TensorFlowImageClassifier` which was initialized before invokes the method `recognizeImage(croppedBitmap)` on the preprocessed image. Within this method, the image data is transferred from `int` to `float`. Then, the previous initialized object of the `TensorflowInferenceInterface` calls the method `feed` to pass the float values to the input layer of the model. Afterwards, the `inferenceInterface` invokes the method `run` with the specified output layer in order to process the classification. Subsequently, the output data is fetched by calling `fetch` on the output layer listing 5.

```

1
2 // Copy the input data into TensorFlow.
3 Trace.beginSection("feed");
4 inferenceInterface.feed(inputName, floatValues, 1, inputSize, inputSize, 3);
5 Trace.endSection();
6
7 // Run the inference call.
8 Trace.beginSection("run");
9 inferenceInterface.run(outputNames, logStats);
10 Trace.endSection();
11
12 // Copy the output Tensor back into the output array.
13 Trace.beginSection("fetch");
14 inferenceInterface.fetch(outputName, outputs);
15 Trace.endSection();

```

Listing 3: Classifying images by the `inferenceInterface`

Afterwards, the results are reordered according to the highest probability and returned to the `results` field for displaying the output.

4.7. Deployment and Validation

Alice

5. Evaluation

Andi - prio von niedrig zu hoch

- regarding implementation time
- regarding performance
- regarding quality in accuracy

- handy performance?

6. Conclusion

- tutorials not complete, different
- which model is better
- prospects, improvements, Recommendations

Beispiele frs referenzieren:

In Figure 9 ist das HS Mnchen Logo zu sehen.



Figure 9: FH-Logo

Oder auch eines Codes wie in listing 6.

```

1
2 bottleneck_path_2_bottleneck_values = {}
3
4
5 def create_bottleneck_file(bottleneck_path, image_lists, label_name, index,
6                             image_dir, category, sess, jpeg_data_tensor,
7                             decoded_image_tensor, resized_input_tensor,
8                             bottleneck_tensor):
9     """Create a single bottleneck file."""
10    tf.logging.info('Creating bottleneck at ' + bottleneck_path)
11    image_path = get_image_path(image_lists, label_name, index,
12                                image_dir, category)
13    if not gfile.Exists(image_path):
14        tf.logging.fatal('File does not exist %s', image_path)
15    image_data = gfile.GFile(image_path, 'rb').read()
16    try:
17        bottleneck_values = run_bottleneck_on_image(
18            sess, image_data, jpeg_data_tensor, decoded_image_tensor,
19            resized_input_tensor, bottleneck_tensor)
20    except Exception as e:
21        raise RuntimeError('Error during processing file %s (%s)' % (image_path,
22                                                                    str(e)))
23    bottleneck_string = ','.join(str(x) for x in bottleneck_values)
24    with open(bottleneck_path, 'w') as bottleneck_file:
25        bottleneck_file.write(bottleneck_string)

```

Listing 4: Some python code

Sectionrefs: In section 2 ist vieles noch nicht fertig.

SubSectionrefs: In section 2.1 wird dann nher auf den Inhalt eingegangen.

SubSubSectionrefs: In section 4.3.2 gehts ans eingemachte.

Beispiele frs zitieren:

Fr einen noch besseren berblick, kann das Buch von ? oder auch andere refs wie ? hinzugezogen werden (?).

Wenn in klammern und Seitenzahl (?, p. 3)

als compared, aber ohne Seitenzahl (cmp. ?)

als compared mit Seitenzahl, das nd heit "no date", da keine Jahreszahl vorhanden (cmp. ?, p. 5)

List of Figures

1.	Model architecture of LeNet-5	3
2.	ImageNet Classification top-5 error	3
3.	Model architecture of AlexNet	4
4.	Model architecture of VGG19	4
5.	Model architecture of Inception	5
6.	Inception modules	5
7.	ResNet block	6
8.	Standard convolution and depthwise seperable convolution	6
9.	FH-Logo	12

A. Anhang

A.1. Anhang vom Anhang