

# Image Classification regarding objects in a home environment for the RoboCup@Home competition

Matteo Facci  
1597454

October 6, 2022

## 1 Introduction

This report provides some solutions for the classification problem of images regarding objects in a home environment.

To this end, the RoboCup@Home-Objects dataset is used, which has been developed within the RoboCup@Home competition. Starting from the manipulation of the latter, in order to prepare the environment where the program will be executed, new methods for training a classification model and predicting new results will be shown and compared.

There is a total amount of 8 classes :

```
classes = ['glass food container', 'potato crisps', 'sports drinking water', 'pasta sides', 'replacement heads',  
'steak knives', 'plums', 'side dish plate'];
```

## 2 Programming Language

The program is developed in python in order to make use of *Keras* Library, a powerful library for developing and evaluating deep learning models. It wraps the efficient numerical computation libraries *Theano* and *TensorFlow* and allows to define and train neural network models in just a few lines of code.

Colab is used as a development environment. It allows to write and execute arbitrary python code through the browser, and it is especially well suited to machine learning and data analysis. More technically, Colab is a hosted Jupyter notebook service that requires no setup to use, while providing free access to computing resources including GPUs.

Moreover a particular section of the report has been developed with the help of Matlab and its Deep Learning Toolbox, which provides a framework for designing and implementing deep neural networks with algorithms, pretrained models, and useful apps.

## 3 Data pre-processing and data augmentation

The Dataset consists on 8779 images divided not equally for each of the 8 categories mentioned above.

That is very few examples to learn from for a classification problem, due to the disproportionate number of classes compared to the insufficient number of images. So this is a

challenging machine learning problem, but it is also a realistic one: in a lot of real-world use cases, data collection can be extremely expensive or sometimes near-impossible.

It can almost be considered an unrepresentative dataset. An unrepresentative dataset means a dataset that may not capture the statistical characteristics relative to another dataset drawn from the same domain, such as between a train and a validation dataset. This can commonly occur if the number of samples in a dataset is too small, relative to another dataset.

In order to make the most of the few training examples, they are "augmented" via a number of random transformations, so that the model would never see twice the exact same picture. This helps prevent overfitting and helps the model generalize better.

### 3.1 Data augmentation in Python

In *Keras* this can be done via the *ImageDataGenerator* class. This class allows to:

- configure random transformations and normalization operations to be done on your image data during training
- instantiate generators of augmented image batches (and their labels) via `flow_from_directory(directory)`.  
These generators can then be used with the *Keras* model methods that accept data generators as inputs, `fit_generator`, `evaluate_generator` and `predict_generator`.

For this purpose the contents of the dataset folder have been divided into two subfolders, the *train* folder, containing the 80% of the total amount of the images in the dataset (7026 training samples from 8 classes), and the *test* folder, containing the remaining 20% (1753 test samples from 8 classes).

The final structure of the folder is illustrated below.

```
train
├── plums
├── glass food container
├── potato crisps
├── sports drinking water
├── pasta sides
├── replacement heads
├── steak knives
├── plums
├── side dish plate
└── test
    ├── plums
    ├── ...
    └── side dish plate
```

All the images are loaded in color mode RGB with batch size 32. The target size of *ImageDataGenerator* is one of the parameters which has been changed in the experiments basing on the input shape of the Convolutional Neural Network in exam. When this parameter is present, the images are stretched to the target size.

### 3.2 Data augmentation in Matlab

The logic in Matlab is practically the same but, thanks to the Deep Learning Toolbox, each operation can be performed with a single function.

`augmentedImageDatastore` transforms batches of training and test data, with optional pre-processing such as resizing, rotation, and reflection. It resizes images to make them compatible with the input size of the network in exam.

In Matlab, the contents of the dataset folder have not been divided into subfolders for *train* and *test*, since the function `splitEachLabel` operates the division automatically. Also in this case, the 80% of the total amount of the images are used for training and the remaining 20% for testing. All the images are loaded in color mode RGB with batch size 32.

## 4 Training a neural network from scratch

First a new neural network from scratch is built for doing classification on the dataset. This will give an initial baseline to compare to other techniques used and to other models of neural networks. For sake of simplicity it has been called *Netfs* (from *Net from scratch*).

One of the simplest neural networks is *LeNet*, but certainly, given its structure, this does not allow to classify in a sufficiently accurate way the type of images contained in the dataset. However, this can be a starting point for the creation of *Netfs*.

This is the difference between the *LeNet* layers and the new *Netfs*' ones:

LeNet	Netfs
Convolutional 6 kernels 5x5	Convolutional 12 kernels 11x11
Average Pooling 2x2 stride 2x2	Average Pooling 2x2 stride 2x2
Convolutional 16 kernels 5x5	Convolutional 32 kernels 8x8
Average Pooling 2x2 stride 2x2	Average Pooling 2x2 stride 2x2
Convolutional 120 kernels 5x5	Convolutional 240 kernels 8x8
Fully connected, 84 units	Fully connected, 2048 units (dropout 0.5)
-	Fully connected, 512 units
Fully connected, 24 units	Fully connected, 8 units

All of the layers in *Netfs* use *ReLU* as activation function except for the last layer which uses *softmax*, as in *LeNet*. The size of kernels has been enlarged in order to test if this can be a tuning that allows the net to perform better. After all layers, *BatchNormalization* and a larger *dense* layer have been added, and this helped to create a deeper network which performed better.

In fact a deeper neural network needs less neurons to reach sufficient results, so it is a way to increase the results quality without increasing the training effort and the convergence time.

The loss function will be *categorical cross – entropy loss*, and the learning algorithm will be *Adam*. Various things about this network can be changed to get better performance, perhaps using a larger network or a different optimizer will help, but for the purposes, the goal is to get an approximate baseline for comparison’s sake.

During the training process, after the creation of the model, in order to prevent overfitting, **ModelCheckpoint** and **EarlyStopping** methods are imported from *Keras*. An object of both is created and passed as callback function to **fit\_generator**. **ModelCheckpoint** allows to save the model by monitoring a specific parameter (validation accuracy in this case).

The model will only be saved if the validation accuracy in the current epoch is greater than the one in the last epoch. **EarlyStopping** instead stops the training of the model earlier if there is no increase in the parameter monitored. The *patience* parameter is set to 15 which means that the model will stop to train if it does not see any rise in validation accuracy in 15 epochs.

The same training approach has been used also for training from scratch two state-of-the-art neural networks for comparison’s purpose.

## 4.1 Netfs

- *input shape* :  $224 \times 224$
- *total params* : 2078872
- *trainable params* : 2073184
- *non – trainable params* : 5688
- *epochs* : 80

Metric	Value
<i>Accuracy</i>	0.651
<i>Precision</i>	0.678
<i>Recall</i>	0.653
<i>F1 Score</i>	0.658

## 4.2 AlexNet

- *input shape* :  $227 \times 227$
- *total params* : 35427792
- *trainable params* : 35406656
- *non – trainable params* : 21136
- *epochs* : 40

Metric	Value
<i>Accuracy</i>	0.610
<i>Precision</i>	0.637
<i>Recall</i>	0.610
<i>F1 Score</i>	0.614

### 4.3 VGG16

- *input shape* :  $224 \times 224$
- *total params* : 134293320
- *trainable params* : 134293320
- *non – trainable params* : 0
- *epochs* : 40

Metric	Value
<i>Accuracy</i>	0.681
<i>Precision</i>	0.693
<i>Recall</i>	0.681
<i>F1 Score</i>	0.684

## 5 Transfer Learning

Another approach to the image classification problem is to fine-tune deeper layers by training the network on the new dataset with the pre-trained network as a starting point.

Fine-tuning a network with transfer learning is often faster and easier than constructing and training a new network from scratch. The network has already learned a rich set of image features, but in this way, it can learn features specific to the new dataset.

This method usually gives the highest accuracy and using a pre-trained network with transfer learning is typically much faster and easier than training a network from scratch.

In the present case a *VGG16* model pre-trained with ImageNet has been used in order to transfer learning to a new neural network called *TransferNet*.

The classification before the transfer occurs in the second part of the *VGG16*'s model, which takes the image features in input and picks a category.

This classifier part contains two hidden *fully connected* (or *dense*) layers, each with 4096 neurons and a dense layer with 1000 neurons, one for each of the ImageNet categories.

The *softmax* activation function is used for these neurons, so that the 1000 values they spit out sum up to unity, and can be considered as probabilities. The transfer learning procedure for *VGG16* and the other pre-trained models, consists on downloading the model with the weights resulting from the training on ImageNet.

Then the classifier part is replaced by the *TransferNet* classifier, adapted to the present problem. This classifier will have 8 output neurons in the last layer, one for each class of the dataset. Finally, all the layers of the convolutional part will be freezed, so that only the parameters of the new classifier will be trained.

## 5.1 ImageNet Database

The majority of the pre-trained networks are trained on a subset of the ImageNet database, which is used in the ImageNet Large-Scale Visual Recognition Challenge (ILSVRC). It features over 14 million images sorted in 1000 object categories.

## 5.2 VGG16-TransferNet

- *input shape* :  $224 \times 224$
- *total params* : 19855348
- *trainable params* : 7449492
- *non – trainable params* : 12405856
- *epochs* : 48

Metric	Value
<i>Accuracy</i>	0.785
<i>Precision</i>	0.804
<i>Recall</i>	0.782
<i>F1 Score</i>	0.788

## 6 Feature Extraction

Feature extraction is an easy and fast way to use the power of deep learning without investing time and effort into training a full network.

It only requires a single pass over the training images, therefore it is especially useful in absence of a GPU.

In this regard, the entire feature extraction operation was carried out in this case using Matlab for comparison, on a laptop with single CPU (Intel(R) Core(TM) i5-8265U CPU @ 1.60GHz - RAM 8GB).

The learned image features are extracted by using a pre-trained network, and then those features are used to train a classifier, such as a support vector machine (SVM) using `fitcecoc`, a useful function for fitting multiclass models included in the Statistics and Machine Learning Toolbox.

Pros of this technique are: a fast training and, if the data is very similar to the original data (ImageNet), then the more specific features extracted deeper in the network are likely to be useful for the new task. About the cons: if the data is very different from the original data (ImageNet), then the features extracted deeper in the network might be less useful for the new task.

Taking into account the following graph, all the pre-trained networks on ImageNet represented there have been used for the feature extraction, and in this case, in addition to the results also the feature extraction and training time has been reported.

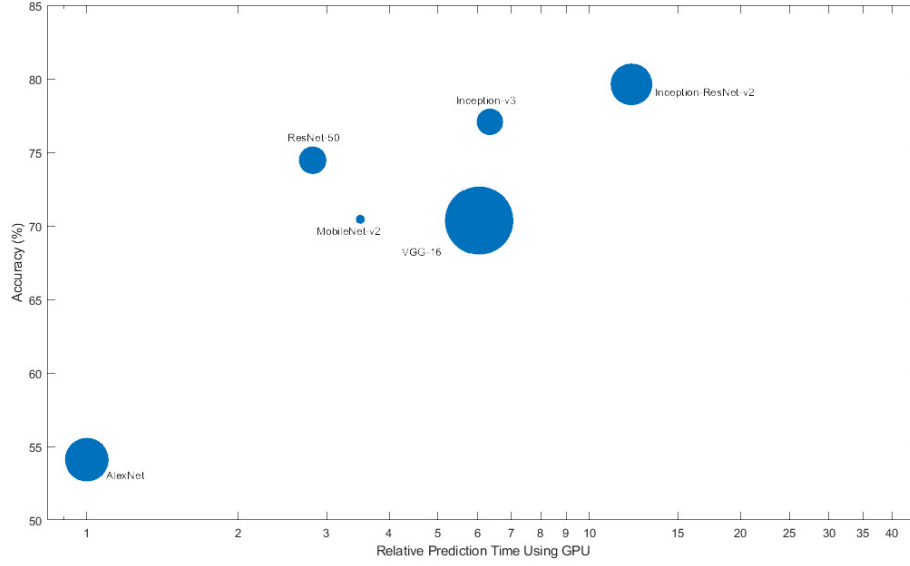


Figure 1: Indication of the relative speeds (using GPU) and accuracy of the different networks trained on the ImageNet dataset. The greater the circle, the greater the weight of the net.

## 6.1 Alexnet

*Training time* = 102.547456 s

Metric	Value
<i>Precision</i>	0.727564
<i>Recall</i>	0.792519
<i>F1 Score</i>	0.724455

		Confusion Matrix									
Output Class	glass food container	158 8.1%	3 0.2%	0 0.0%	2 0.1%	5 0.3%	2 0.1%	4 0.2%	0 0.0%	90.8%	9.2%
	pasta sides	0 0.0%	168 8.6%	0 0.0%	17 0.9%	7 0.4%	3 0.2%	4 0.2%	0 0.0%	84.4%	15.6%
	plums	2 0.1%	5 0.3%	229 11.7%	6 0.3%	3 0.2%	9 0.5%	1 0.1%	0 0.0%	89.8%	10.2%
	potato crisps	6 0.3%	15 0.8%	4 0.2%	181 9.3%	3 0.2%	20 1.0%	1 0.1%	4 0.2%	77.4%	22.6%
	replacement heads	5 0.3%	2 0.1%	0 0.0%	0 0.0%	145 7.4%	7 0.4%	11 0.6%	5 0.3%	82.9%	17.1%
	side dish plate	63 3.2%	39 2.0%	9 0.5%	31 1.6%	59 3.0%	199 10.2%	17 0.9%	22 1.1%	45.3%	54.7%
	sports drinking water	9 0.5%	11 0.6%	2 0.1%	7 0.4%	14 0.7%	2 0.1%	205 10.5%	2 0.1%	81.3%	18.7%
	steak knives	1 0.1%	1 0.1%	0 0.0%	0 0.0%	8 0.4%	2 0.1%	1 0.1%	211 10.8%	94.2%	5.8%
		64.8%	68.9%	93.9%	74.2%	59.4%	81.6%	84.0%	86.5%	76.6%	
		35.2%	31.1%	6.1%	25.8%	40.6%	18.4%	16.0%	13.5%	23.4%	
		Target Class									
		glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives		

Figure 2: *Confusion Matrix* for SVM using pre-trained AlexNet.

## 6.2 VGG16

*Training time* = 1800.486780 s

Metric	Value
<i>Precision</i>	0.805769
<i>Recall</i>	0.811865
<i>F1 Score</i>	0.806679

**Confusion Matrix**

Output Class	glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives	
glass food container	172 11.0%	3 0.2%	3 0.2%	2 0.1%	3 0.2%	19 1.2%	1 0.1%	2 0.1%	83.9%
pasta sides	7 0.4%	137 8.8%	6 0.4%	33 2.1%	11 0.7%	18 1.2%	0 0.0%	12 0.8%	61.2%
plums	1 0.1%	2 0.1%	174 11.2%	5 0.3%	0 0.0%	6 0.4%	3 0.2%	0 0.0%	91.1%
potato crisps	1 0.1%	19 1.2%	2 0.1%	135 8.7%	1 0.1%	6 0.4%	3 0.2%	2 0.1%	79.9%
replacement heads	2 0.1%	10 0.6%	1 0.1%	3 0.2%	158 10.1%	6 0.4%	4 0.3%	11 0.7%	81.0%
side dish plate	8 0.5%	18 1.2%	7 0.4%	17 1.1%	2 0.1%	132 8.5%	1 0.1%	2 0.1%	70.6%
sports drinking water	4 0.3%	6 0.4%	2 0.1%	0 0.0%	19 1.2%	8 0.5%	183 11.7%	0 0.0%	82.4%
steak knives	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	0 0.0%	166 10.6%	99.4%
	88.2%	70.3%	89.2%	69.2%	81.0%	67.7%	93.8%	85.1%	80.6%
	11.8%	29.7%	10.8%	30.8%	19.0%	32.3%	6.2%	14.9%	19.4%
	glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives	
	Target Class								

Figure 3: *Confusion Matrix* for SVM using pre-trained VGG16.

## 6.3 MobileNet-v2

*Training time* = 331.062729 s

Metric	Value
<i>Precision</i>	0.823770
<i>Recall</i>	0.826767
<i>F1 Score</i>	0.823620



		Confusion Matrix									
Output Class	glass food container	217 11.1%	6 0.3%	0 0.0%	4 0.2%	2 0.1%	14 0.7%	5 0.3%	3 0.2%	86.5% 13.5%	
	pasta sides	4 0.2%	158 8.1%	6 0.3%	13 0.7%	8 0.4%	9 0.5%	0 0.0%	2 0.1%	79.0% 21.0%	
	plums	3 0.2%	3 0.2%	223 11.4%	10 0.5%	3 0.2%	4 0.2%	3 0.2%	0 0.0%	89.6% 10.4%	
	potato crisps	2 0.1%	35 1.8%	3 0.2%	177 9.1%	5 0.3%	13 0.7%	1 0.1%	0 0.0%	75.0% 25.0%	
	replacement heads	2 0.1%	5 0.3%	3 0.2%	1 0.1%	201 10.3%	8 0.4%	9 0.5%	7 0.4%	85.2% 14.8%	
	side dish plate	13 0.7%	33 1.7%	7 0.4%	34 1.7%	4 0.4%	185 9.5%	4 0.2%	5 0.3%	64.2% 35.8%	
	sports drinking water	2 0.1%	4 0.2%	2 0.1%	5 0.3%	12 0.6%	6 0.3%	222 11.4%	2 0.1%	87.1% 12.9%	
	steak knives	1 0.1%	0 0.0%	0 0.0%	0 0.0%	6 0.3%	5 0.3%	0 0.0%	225 11.5%	94.9% 5.1%	
		88.9% 11.1%	64.8% 35.2%	91.4% 8.6%	72.5% 27.5%	82.4% 17.6%	75.8% 24.2%	91.0% 9.0%	92.2% 7.8%	82.4% 17.6%	
		Target Class									
	glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives			

Figure 4: *Confusion Matrix* for SVM using pre-trained MobileNet-v2.

## 6.4 ResNet-50

*Training time* = 633.149611 s

Metric	Value
<i>Precision</i>	0.866667
<i>Recall</i>	0.875311
<i>F1 Score</i>	0.868155

		Confusion Matrix									
Output Class	glass food container	229 11.7%	10 0.5%	2 0.1%	1 0.1%	3 0.2%	17 0.9%	5 0.3%	3 0.2%	84.8% 15.2%	
	pasta sides	1 0.1%	149 7.6%	1 0.1%	7 0.4%	6 0.3%	10 0.5%	0 0.0%	1 0.1%	85.1% 14.9%	
	plums	0 0.0%	1 0.1%	226 11.6%	1 0.1%	0 0.0%	3 0.2%	2 0.1%	0 0.0%	97.0% 3.0%	
	potato crisps	0 0.0%	56 2.9%	12 0.6%	218 11.2%	9 0.5%	27 1.4%	5 0.3%	2 0.1%	66.3% 33.7%	
	replacement heads	0 0.0%	2 0.1%	0 0.0%	1 0.1%	182 9.3%	5 0.3%	2 0.1%	4 0.2%	92.9% 7.1%	
	side dish plate	6 0.3%	15 0.8%	2 0.1%	12 0.6%	8 0.4%	175 9.0%	0 0.0%	7 0.4%	77.8% 22.2%	
	sports drinking water	7 0.4%	9 0.5%	0 0.0%	3 0.2%	30 1.5%	7 0.4%	230 11.8%	1 0.1%	80.1% 19.9%	
	steak knives	1 0.1%	2 0.1%	1 0.1%	1 0.1%	6 0.3%	0 0.0%	0 0.0%	226 11.6%	95.4% 4.6%	
		93.9% 6.1%	61.1% 38.9%	92.6% 7.4%	89.3% 10.7%	74.6% 25.4%	71.7% 28.3%	94.3% 5.7%	92.6% 7.4%	83.8% 16.2%	
			Target Class								
		glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives		

Figure 5: *Confusion Matrix* for SVM using pre-trained ResNet-50.

## 6.5 Inception-v3

*Training time* = 1256.542420 s

Metric	Value
<i>Precision</i>	0.836066
<i>Recall</i>	0.838652
<i>F1 Score</i>	0.835127

**Confusion Matrix**

Output Class	glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives	
glass food container	231 11.8%	6 0.3%	4 0.2%	6 0.3%	5 0.3%	10 0.5%	7 0.4%	2 0.1%	85.2% 14.8%
pasta sides	2 0.1%	176 9.0%	3 0.2%	19 1.0%	9 0.5%	17 0.9%	4 0.2%	2 0.1%	75.9% 24.1%
plums	2 0.1%	6 0.3%	217 11.1%	14 0.7%	0 0.0%	2 0.1%	5 0.3%	2 0.1%	87.5% 12.5%
potato crisps	1 0.1%	16 0.8%	3 0.2%	162 8.3%	2 0.1%	7 0.4%	1 0.1%	0 0.0%	84.4% 15.6%
replacement heads	1 0.1%	4 0.2%	1 0.1%	6 0.3%	210 10.8%	9 0.5%	11 0.6%	9 0.5%	83.7% 16.3%
side dish plate	5 0.3%	26 1.3%	9 0.5%	33 1.7%	9 0.5%	195 10.0%	1 0.1%	3 0.2%	69.4% 30.6%
sports drinking water	2 0.1%	9 0.5%	6 0.3%	1 0.1%	4 0.2%	0 0.0%	215 11.0%	0 0.0%	90.7% 9.3%
steak knives	0 0.0%	1 0.1%	1 0.1%	3 0.2%	5 0.3%	4 0.2%	0 0.0%	226 11.6%	94.2% 5.8%
	94.7% 5.3%	72.1% 27.9%	88.9% 11.1%	66.4% 33.6%	86.1% 13.9%	79.9% 20.1%	88.1% 11.9%	92.6% 7.4%	83.6% 16.4%
	glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives	
	Target Class								

Figure 6: *Confusion Matrix* for SVM using pre-trained Inception-v3.

## 6.6 Inception-Resnet-v2

*Training time* = 2151.029954 s

Metric	Value
<i>Precision</i>	0.842213
<i>Recall</i>	0.853586
<i>F1 Score</i>	0.842395

		Confusion Matrix									
Output Class	glass food container	220 11.3%	8 0.4%	3 0.2%	2 0.1%	0 0.0%	9 0.5%	4 0.2%	1 0.1%	89.1% 10.9%	
	pasta sides	3 0.2%	157 8.0%	1 0.1%	9 0.5%	0 0.0%	1 0.0%	1 0.1%	0 0.0%	91.3% 8.7%	
	plums	1 0.1%	2 0.1%	231 11.8%	9 0.5%	0 0.0%	2 0.1%	4 0.2%	0 0.0%	92.8% 7.2%	
	potato crisps	0 0.0%	30 1.5%	1 0.1%	186 9.5%	0 0.0%	20 1.0%	1 0.1%	2 0.1%	77.5% 22.5%	
	replacement heads	3 0.2%	21 1.1%	2 0.1%	11 0.6%	11 11.7%	13 0.7%	14 0.7%	14 0.7%	74.6% 25.4%	
	side dish plate	13 0.7%	22 1.1%	6 0.3%	23 1.2%	7 0.4%	196 10.0%	3 0.2%	18 0.9%	68.1% 31.9%	
	sports drinking water	4 0.2%	4 0.2%	0 0.0%	4 0.2%	3 0.2%	0 0.0%	217 11.1%	1 0.1%	91.9% 8.1%	
	steak knives	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	3 0.2%	0 0.0%	208 10.7%	97.7% 2.3%	
		90.2% 9.8%	64.3% 35.7%	94.7% 5.3%	76.2% 23.8%	93.9% 6.1%	80.3% 19.7%	88.9% 11.1%	85.2% 14.8%	84.2% 15.8%	
			glass food container	pasta sides	plums	potato crisps	replacement heads	side dish plate	sports drinking water	steak knives	
		Target Class									

Figure 7: *Confusion Matrix* for SVM using pre-trained Inception-Resnet-v2.

## 7 Conclusions

Collected all the data resulted from the usage of different techniques and methods for solving the present image classification problem it is possible to formulate some considerations.

Looking separately at each of the three methods used, the first analyzed is to train a neural network from scratch.

This is the method which provided the worst results among the three, in fact, due probably to the quite unrepresentative dataset and to the consequent risk of overfitting, the best result achieved was the one of the state-of-the-art *VGG16*, with a *F1 Score* of 0.684, slightly above the one obtained by *Netfs*. Poor *F1 Score* of 0.614 for *AlexNet*.

Not exactly satisfactory results accompanied by a training time exceeding 2 hours for all three models analyzed.

By using a pre-trained *VGG16* model on ImageNet and performing the Transfer Learning technique with the *TransferNet* neural network a 10% higher score is obtained, in particular an *F1 Score* of 0.788, and a certainly reduced time but always much higher than in the case of the Feature Extraction.

In this case, looking at the confusion matrices and at the *F1 Score* of each model, it is evident that Support Vector Machine models perform the best (for the given problem, with the given dataset) in terms of correct predictions, since their scores are of 0.868155 for the *ResNet – 50*-based model, the highest score obtained among the ones of all the analyzed classifiers, and of 0.724455 for the worst case given also in this case by *AlexNet*.

As anticipated, the feature extraction and training time for the case of the best model, the *ResNet – 50*-based model, is 633.149611 s, extremely reduced if compared to the times obtained with previous techniques. And considering also the limited hardware used, the technique and the model in question are among the best alternatives to solve the problem in the shortest possible time, and with the least waste of resources. Thanks also to the ease of use of Matlab and its prepared toolboxes.