

CNN Classifier

Giovanni Battilana, Marzia Paschini, Marco Sgobino

February 2, 2022

Contents

1	Assignment general description	2
2	Goal of this paper	2
3	Adopted tools	2
4	Training a baseline network	2
4.1	Importing the dataset	3
4.2	Requirements	3
4.3	Splitting the dataset	4
4.4	Resizing images	4
4.5	Layers definition	4
4.6	Network training options	5
4.7	Training the network and obtaining the accuracy	6
4.8	Choice of MaxEpoch and stopping criterion	6
4.9	Fine-tuning initial learning rate	6
4.10	Baseline training progress and confusion matrix	7
5	Improving the network	10
5.1	Data augmentation	10
5.2	Batch normalization	10
5.3	Improving convolutional layers	11
5.4	Modifying network layout	11
5.4.1	Adding fully-connected and dropout layers	11
5.4.2	MaxEpoch	11
5.4.3	Results	12
5.5	Parameters and layout enhancements that did not provide benefit	12
5.6	Comparison with the baseline	13
5.7	Ensemble network based on previous improvements	16
5.8	Justifying the added complexity	16
6	Adopting transfer learning with AlexNet	17
6.1	Transfer learning by freezing weights	17

6.1.1	Input layer	18
6.1.2	Training options	19
6.1.3	Results	19
6.2	AlexNet as feature extractor	22
7	Summary of results	24
	References	25

1 Assignment general description

This project requires the implementation of an image classifier based on Convolutional Neural Networks. The provided dataset[1], contains 15 categories (office, kitchen, living room, bedroom, store, industrial, tall building, inside city, street, highway, coast, open country, mountain, forest, suburb), and is already divided in training set and test set.

2 Goal of this paper

This paper had 3 goals of increasing difficulty:

1. building a *shallow network* with a layout specified from the Teacher in order to have a *baseline* network to allow further comparisons. Such baseline should be the reference from which subsequent improvements are evaluated;
2. improving the previous baseline network and compare results with the baseline. Comments on strategies adopted should be provided;
3. adopting *transfer learning* based on the pre-trained network *AlexNet*. Comments on performance improvements should be given.

3 Adopted tools

Students adopted the programming language and numeric computing environment **MATLAB** as both text editor and simulation software to describe, implement and build the trained networks.

Specific MATLAB Toolbox were needed – the project should require *Deep Learning Toolbox* (along with *Deep Learning Toolbox Model for AlexNet Network*), *Computer Vision Toolbox*, *Image Processing Toolbox*, and *Statistics and Machine Learning Toolbox*. Optionally, to improve processing speed during the Convolutional Neural Network build phase, students installed and enabled the *Parallel Computing Toolbox*.

4 Training a baseline network

The first step was the building of a baseline network (a shallow network), whose layout had been assigned by the Teacher. In practice, this phase required to strictly follow given requirements in order to obtain *an overall test accuracy of around 30%*.

4.1 Importing the dataset

The training dataset comprises 1500 images, of 15 categories (office, kitchen, living room, bedroom, store, industrial, tall building, inside city, street, highway, coast, open country, mountain, forest, suburb). Each photo lies inside a directory, whose name denotes the category. Similarly, the test dataset comprises 2985 pictures lying in the same categories as above.

In order to import the dataset into an *image datastore* object, students used the following instructions,

```
TrainDatasetPath = fullfile('dataset', 'train');  
imds = imageDatastore(TrainDatasetPath, 'IncludeSubfolders', true, 'LabelSource',  
    'foldernames');
```

Each image in the dataset had its own – non common – size. Moreover, images do not share a common aspect ratio. In order to obtain a shared size and aspect ratio, the group has chosen to follow the simple approach of *rescaling the whole image independently along the vertical and the horizontal axis*, in order to achieve the proper size. As will be shown later, the chosen aspect ratio will be of 1 : 1 (a square).

4.2 Requirements

Layout of the network is described in the following Table 1:

#	Layer type	Size and parameters
1	Image input	$64 \times 64 \times 1$ images
2	Convolution	8 3×3 convolutions, stride 1
3	ReLU	
4	Max Pooling	2×2 max pooling, stride 2
5	Convolution	16 3×3 convolutions, stride 1
6	ReLU	
7	Max Pooling	2×2 max pooling, stride 2
8	Convolution	32 3×3 convolutions, stride 1
9	ReLU	
10	Fully Connected	15
11	Softmax	softmax
12	Classification	crossentropyex

Table 1: Baseline Layout

Other non-optional requirements were the following ones,

- using input images of size 64×64 (as shown in layer 1 of Table 1);
- use 85% of the provided dataset for the training set and the remaining portion for the validation set, so that each label is properly split with the same percentage;
- employ **stochastic gradient descent with momentum** as the optimization algorithm;
- adopt *minibatches* of size 32;

- initial weights should be drawn from a Gaussian distribution with a mean of 0 and a standard deviation of 0.01. Initial bias values should be set to 0;
- choose a stopping criterion.

4.3 Splitting the dataset

To split the dataset as requested, students wrote and ran the following code:

```
trainQuota=0.85;
[imdsTrain, imdsValidation] = splitEachLabel(imds, trainQuota, 'randomize');
```

These instructions were sufficient to obtain two, distinct, datasets – the first one for training, the second one for validation. Images should be elected and put in the datasets according to a random process.

4.4 Resizing images

In order to resize images to match the requested size of 64×64 , students ran:

```
imds.ReadFcn = @(x)imresize(imread(x), [64 64]);
```

This command should overload the ReadFcn function in imds, so that instead of simply reading images they should be resized too.

4.5 Layers definition

Layers layout is specified in the provided Table 1. In order to obtain such layout, the students organized the network layout information in a single object named layers, such as

```
layers = [
    imageInputLayer([64 64 1], 'Name', 'input')

    convolution2dLayer(3,8, 'Padding','same', 'Stride', [1 1], 'Name','conv_1',...
        'WeightsInitializer', @(sz) randn(sz)*0.01,...
        'BiasInitializer', @(sz) zeros(sz))

    reluLayer('Name','relu_1')

    maxPooling2dLayer(2, 'Stride',2,...
        'Name','maxpool_1')

    convolution2dLayer(3,16, 'Padding','same', 'Stride', [1 1], 'Name','conv_2',...
        'WeightsInitializer', @(sz) randn(sz)*0.01,...
        'BiasInitializer', @(sz) zeros(sz))

    reluLayer('Name','relu_2')

    maxPooling2dLayer(2, 'Stride',2,...
```

```

    'Name', 'maxpool_2')

convolution2dLayer(3,32, 'Padding','same', 'Stride', [1 1],...
    'Name','conv_3',...
    'WeightsInitializer', @(sz) randn(sz)*0.01,...
    'BiasInitializer', @(sz) zeros(sz))

reluLayer('Name','relu_3')

fullyConnectedLayer(15, 'Name','fc_1',...
    'WeightsInitializer', @(sz) randn(sz)*0.01,...
    'BiasInitializer', @(sz) zeros(sz))

softmaxLayer('Name','softmax')
classificationLayer('Name','output')
];

```

A peculiar aspect of this representation is the weights initialization, which is handled by the function `handle @(sz) randn(sz)*0.01`, which should yield proper random numbers for normally-distributed weights initialization with standard deviation 0.01.

Similarly, `@(sz) zeros(sz)` will generate the null vectors required to initialize the biases.

A name should be assigned to each layer, while the defined network may be visually inspected with the commands

```

lgraph = layerGraph(layers);
analyzeNetwork(lgraph)

```

4.6 Network training options

Options for training the network should be provided. In order to satisfy the requirements, students set the following options in an object,

```

options = trainingOptions('sgdm', ...
    'InitialLearnRate', InitialLearningRate, ...
    'ValidationData', imdsValidation, ...
    'MiniBatchSize', 32, ...
    'ExecutionEnvironment', 'parallel', ...
    'Plots', 'training-progress' ...
);

```

where the quantity `InitialLearningRate` should be properly fine-tuned.

Each option is necessary, with the sole exception of `ExecutionEnvironment`, which enables parallel computing capabilities:

- `sgdm`: the optimization algorithm in use, as in requirements;

- **ValidationData:** specifies the image datastore to use when validating the error, during the training;
- **MiniBatchSize:** set to 32 as in requirements;
- **Plots:** enables a visual inspection of the training process, useful to spot possible signs of overfitting;

Every other option is left as default.

4.7 Training the network and obtaining the accuracy

The following command will begin the network training with specified layout and options,

```
net = trainNetwork(imdsTrain, layers, options);
```

In order to collect the accuracy with the provided test set, the students implemented the following code,

```
TestDatasetPath = fullfile('dataset', 'test');
imdsTest = imageDatastore(TestDatasetPath, ...
    'IncludeSubfolders', true, ...
    'LabelSource', 'foldernames');
imdsTest.ReadFcn = @(x)imresize(imread(x), [64 64]);

YPredicted = classify(net, imdsTest);
YTest = imdsTest.Labels;

accuracy = sum(YPredicted == YTest) / numel(YTest);
```

where a test dataset should have been created with the same size and characteristics of the training set. The last portion of the snippet realized the prediction phase and the accuracy evaluation.

4.8 Choice of MaxEpoch and stopping criterion

Training the network with an initial learning rate of 0.001 and leaving the default number of epochs (30) unaltered led to evident overfitting of the network, as shown in Figure 1.

With the goal of avoiding overfitting, a lower MaxEpoch value (for instance 8) should be adopted. The value has been obtained by repeatedly testing the network with different values.

4.9 Fine-tuning initial learning rate

Fine-tune the initial learning rate required manual trials along with manual inspection. Values of 0.1, 0.01, 0.001 and 0.0001 (representatives of various order of magnitude) were tried. Only values in the neighborhood of 0.001 gave the required performance of around 30%.

The students decided to elect the network with initial learning rate of 0.001 as the *baseline*, with parameters as in following Table 2:

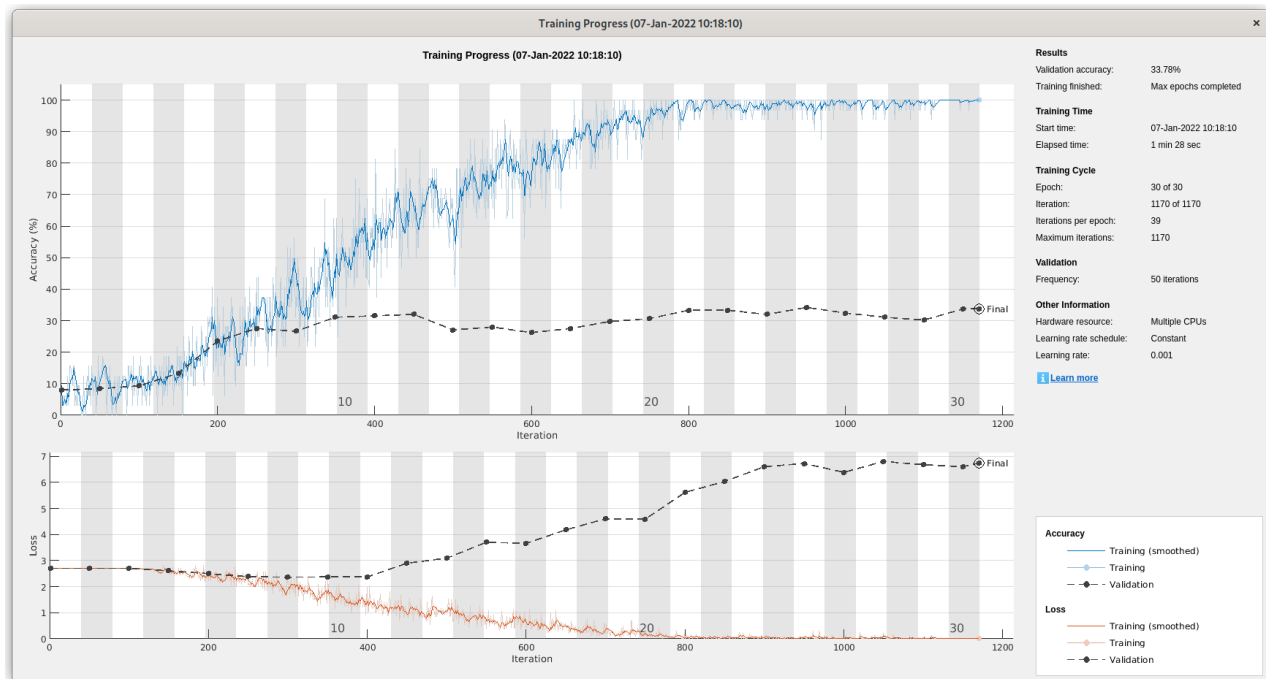


Figure 1: Overfitting of the network when MaxEpoch is the default value of 30. Overfitting is clearly visible after epoch 8, and becomes very relevant after epoch 10.

Parameter	Value
Optimization Algorithm	sgdm
Epochs (stopping criterion)	8
Initial Learning Rate	0.001
MiniBatch size	32
Weights initialization	Gaussian with $\mu = 0$ and $\sigma = 0.01$
Bias initialization	0

Table 2: Baseline parameters.

4.10 Baseline training progress and confusion matrix

The overall training progress for the baseline network is shown in Figure 2.

To obtain the *confusion matrix* for the chosen baseline network, the students ran the following code to first train the network, and then plot a confusion matrix as shown in Figure 3,

```
InitialLearningRate = 0.001;
options = trainingOptions('sgdm', ...
    'InitialLearnRate', InitialLearningRate, ...
    'ValidationData', imdsValidation, ...
    'MiniBatchSize', 32, ...
    'MaxEpochs', 8, ...
    'ExecutionEnvironment', 'parallel', ...
    'Plots', 'training-progress' ...
);

net = trainNetwork(imdsTrain, layers, options);

TestDatasetPath = fullfile('dataset', 'test');
```

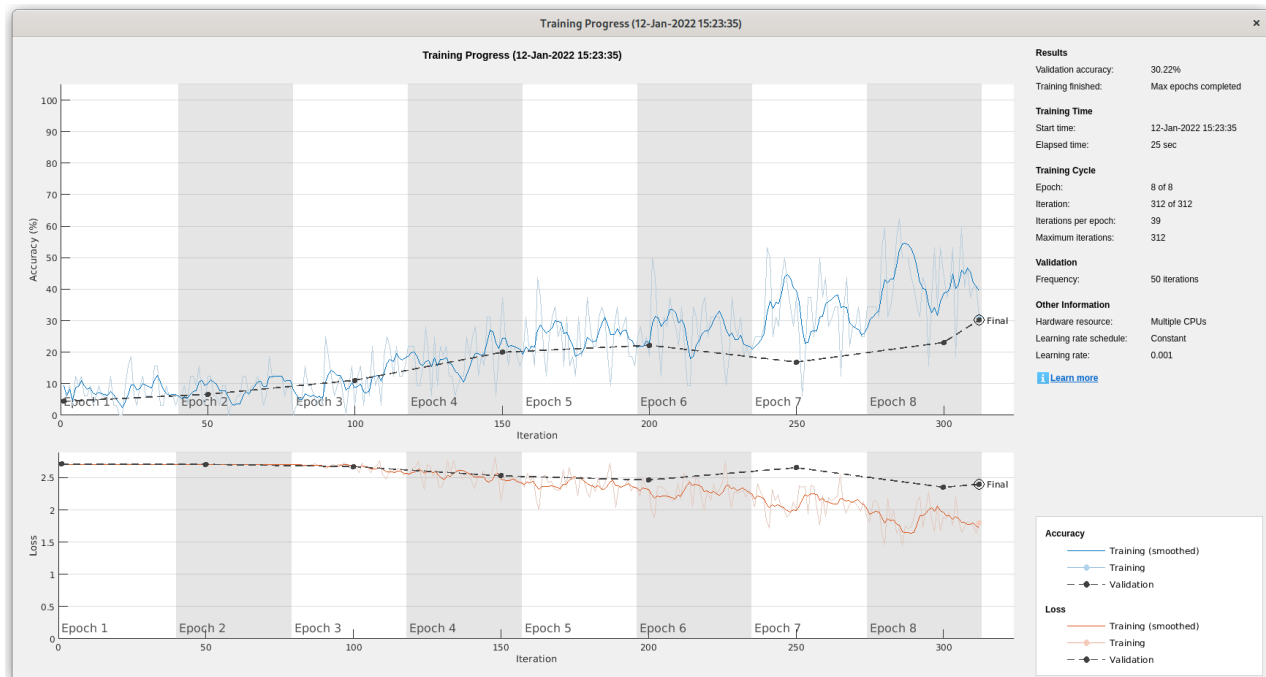


Figure 2: Training progress for the baseline network.

```

imdsTest = imageDatastore(TestDatasetPath, ...
    'IncludeSubfolders',true,'LabelSource','foldernames');
imdsTest.ReadFcn = @(x)imresize(imread(x),[64 64]);

YPredicted = classify(net,imdsTest);
YTest = imdsTest.Labels;

accuracy = sum(YPredicted == YTest)/numel(YTest);

figure
plotconfusion(YTest,YPredicted)

```

The confusion matrix shows that there are some classes that are more often misclassified (*Kitchen, Industrial, Bedroom, Forest, LivingRoom*) and some classes that instead are less often misclassified, such as *Coast, Highway, Office, Street* and *Suburb*. Overall accuracy stands at the value 30.6%: from this point of view, the baseline should represent an improvement from a completely random classifier, whose performance with respect to this classification task is less than a 7%.

The overall code of the baseline network can be read in `baseline_network.m` file.

		Confusion Matrix															
Output Class	Bedroom	5 0.2%	4 0.1%	6 0.2%	0 0.0%	5 0.2%	1 0.0%	2 0.1%	12 0.4%	10 0.3%	1 0.0%	1 0.0%	6 0.2%	2 0.1%	3 0.1%	7 0.2%	7.7% 92.3%
	Coast	25 0.8%	149 5.0%	6 0.2%	39 1.3%	28 0.9%	5 0.2%	5 0.2%	9 0.3%	63 2.1%	12 0.4%	101 3.4%	5 0.2%	5 0.2%	3 0.1%	5 0.2%	32.4% 67.6%
	Forest	1 0.0%	6 0.2%	13 0.4%	0 0.0%	2 0.1%	1 0.0%	2 0.1%	1 0.0%	6 0.2%	2 0.1%	3 0.1%	0 0.0%	1 0.0%	1 0.0%	5 0.2%	29.5% 70.5%
	Highway	1 0.0%	17 0.6%	1 0.0%	89 3.0%	14 0.5%	3 0.1%	3 0.1%	8 0.3%	8 0.3%	1 0.0%	17 0.6%	0 0.0%	4 0.1%	4 0.1%	3 0.1%	51.4% 48.6%
	Industrial	10 0.3%	0 0.0%	5 0.2%	1 0.0%	33 1.1%	8 0.3%	6 0.2%	12 0.4%	11 0.4%	7 0.2%	1 0.0%	3 0.1%	11 0.4%	4 0.1%	12 0.4%	26.6% 73.4%
	InsideCity	6 0.2%	6 0.2%	26 0.9%	1 0.0%	21 0.7%	50 1.7%	4 0.1%	15 0.5%	22 0.7%	4 0.1%	19 0.6%	16 0.5%	12 0.4%	12 0.4%	18 0.6%	21.6% 78.4%
	Kitchen	26 0.9%	6 0.2%	30 1.0%	1 0.0%	19 0.6%	29 1.0%	34 1.1%	40 1.3%	33 1.1%	7 0.2%	12 0.4%	48 1.6%	8 0.3%	10 0.3%	32 1.1%	10.1% 89.9%
	LivingRoom	11 0.4%	0 0.0%	3 0.1%	0 0.0%	9 0.3%	9 0.3%	7 0.2%	23 0.8%	2 0.1%	6 0.2%	1 0.0%	10 0.3%	4 0.1%	1 0.0%	7 0.2%	24.7% 75.3%
	Mountain	2 0.1%	12 0.4%	1 0.0%	1 0.0%	3 0.1%	2 0.1%	0 0.0%	1 0.0%	18 0.6%	2 0.1%	10 0.3%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	34.0% 66.0%
	Office	1 0.0%	14 0.5%	55 1.8%	5 0.2%	16 0.5%	20 0.7%	8 0.3%	10 0.3%	29 1.0%	55 1.8%	13 0.4%	13 0.4%	7 0.2%	2 0.1%	54 1.8%	18.2% 81.8%
	OpenCountry	0 0.0%	26 0.9%	3 0.1%	15 0.5%	9 0.3%	7 0.2%	2 0.1%	0 0.0%	27 0.9%	0 0.0%	94 3.1%	0 0.0%	2 0.1%	1 0.0%	2 0.1%	50.0% 50.0%
	Store	9 0.3%	1 0.0%	34 1.1%	3 0.1%	23 0.8%	36 1.2%	16 0.5%	28 0.9%	9 0.3%	2 0.1%	10 0.3%	82 2.7%	13 0.4%	15 0.5%	22 0.7%	27.1% 72.9%
	Street	4 0.1%	3 0.1%	11 0.4%	3 0.1%	11 0.4%	11 0.4%	6 0.2%	8 0.3%	9 0.3%	3 0.1%	10 0.3%	3 0.1%	109 3.7%	7 0.2%	4 0.1%	54.0% 46.0%
	Suburb	9 0.3%	14 0.5%	13 0.4%	2 0.1%	6 0.2%	18 0.6%	10 0.3%	15 0.5%	16 0.5%	2 0.1%	16 0.5%	19 0.6%	13 0.4%	77 2.6%	3 0.1%	33.0% 67.0%
	TallBuilding	6 0.2%	2 0.1%	21 0.7%	0 0.0%	12 0.4%	8 0.3%	5 0.2%	7 0.2%	11 0.4%	11 0.4%	2 0.1%	10 0.3%	0 0.0%	1 0.0%	82 2.7%	46.1% 53.9%
		4.3% 95.7%	57.3% 42.7%	5.7% 94.3%	55.6% 44.4%	15.6% 84.4%	24.0% 76.0%	30.9% 69.1%	12.2% 87.8%	6.6% 93.4%	47.8% 52.2%	30.3% 69.7%	38.1% 61.9%	56.8% 43.2%	54.6% 45.4%	32.0% 68.0%	30.6% 69.4%
		Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding	
		Target Class															

Figure 3: Confusion matrix for the baseline network.

5 Improving the network

The students propose (and have implemented) various methods in order to almost double the performance of the Convolutional Neural Network from the baseline $\sim 30\%$ to a $\sim 60\%$:

1. *data augmentation*;
2. addition of *batch normalization layers*;
3. changing the *size of the convolutional filters*;
4. modifying network layout.

Individual contribution and net improvement of each of these methods is not investigated – what is provided, instead, is a description of each technique’s implementation along with the necessary details, a brief description of the increasing improvements *at each new addition* and a final overview of all improvements altogether with respect to the baseline.

5.1 Data augmentation

The first improvement of the network is related to *data augmentation*. The students implemented *left-to-right reflections* of the train set, leaving unaltered the test set required for evaluation.

Running code

```
aug = imageDataAugmenter("RandXReflection",true);
imageSize = [64 64 1];
auimds = augmentedImageDatastore(imageSize,imdsTrain,'DataAugmentation',aug);
aunet = trainNetwork(auimds,layers,options);
```

with the same options as the baseline, but MaxEpoch set to 15 (this kind of training resulted to be slower), returned a validation accuracy of around $\sim 39\%$. Hence, the added complexity represented a remarkable improvement from the previous step.

5.2 Batch normalization

The second improvement should consist in adding 3 **Batch Normalization layers** before the ReLU layers. Each layer should be added into layers object, by adding line

```
batchNormalizationLayer('Name','BN_1')
```

before any ReLU layer. The name of the new layers should be changed accordingly to the position of those layers.

The added complexity resulted in a faster learning (a lower MaxEpoch number should be set, from 15 to 6) and a much higher accuracy (51.11% on validation set, 48.7% on test data). With this second addition, the improvements were remarkable too.

5.3 Improving convolutional layers

Convolution layers may be improved by increasing the size of the filter as they are placed towards the output. In practice, this means modifying the convolution layers such that their size increases from 3×3 only, to 3×3 , 5×5 and 7×7 :

```
convolution2dLayer(3,8,'Padding','same','Stride',[1 1], 'Name','conv_1',...
    'WeightsInitializer', @(sz) randn(sz)*0.01,...
    'BiasInitializer', @(sz) zeros(sz))

[...]

convolution2dLayer(5,16,'Padding','same','Stride',[1 1], 'Name','conv_2',...
    'WeightsInitializer', @(sz) randn(sz)*0.01,...
    'BiasInitializer', @(sz) zeros(sz))

[...]

convolution2dLayer(7,32,'Padding','same','Stride',[1 1], 'Name','conv_3',...
    'WeightsInitializer', @(sz) randn(sz)*0.01,...
    'BiasInitializer', @(sz) zeros(sz))
```

This third addition brought a slight improvement, as the validation accuracy is at 55.11%, with the overall test accuracy standing at 50%.

5.4 Modifying network layout

Editing the layout of the network should improve the overall accuracy.

Since layout of the network changed, *MaxEpoch* should be set accordingly, with the goal of avoiding overfitting while leaving the network room for enough complexity.

5.4.1 Adding fully-connected and dropout layers

A new fully-connected layer, along with a ReLU layer, were added:

```
fullyConnectedLayer(256,'Name','fc_1',...
    'WeightsInitializer', @(sz) randn(sz)*0.01,...
    'BiasInitializer', @(sz) zeros(sz))
reluLayer('Name','relu_4')
```

A dropout layer separating the two fully-connected layers was added too,

```
dropoutLayer(.25, 'Name', 'dropout_1')
```

5.4.2 MaxEpoch

In order to let the network have enough room for complexity, a *MaxEpoch* value of 25 should be set.

5.4.3 Results

Overall network layout is illustrated in the following Table 3:

#	Layer type	Size and parameters
1	Image input	$64 \times 64 \times 1$ images
2	Convolution	8 3×3 convolutions, stride 1
3	Batch Normalization	
4	ReLU	
5	Max Pooling	2×2 max pooling, stride 2
6	Convolution	16 5×5 convolutions, stride 1
7	Batch Normalization	
8	ReLU	
9	Max Pooling	2×2 max pooling, stride 2
10	Convolution	32 7×7 convolutions, stride 1
11	Batch Normalization	
12	ReLU	
13	Max Pooling	2×2 max pooling, stride 2
14	Fully Connected	256
15	ReLU	
16	Dropout	25%
17	Fully Connected	15
18	Softmax	softmax
19	Classification	crossentropyex

Table 3: Improved Network layout.

while training options are listed in the table below,

Parameter	Value
Optimization Algorithm	sgdm
Epochs (stopping criterion)	25
Initial Learning Rate	0.001
MiniBatch size	32
Weights initialization	Gaussian with $\mu = 0$ and $\sigma = 0.01$
Bias initialization	0

Table 4: Improved Network parameters.

Students reported that these further improvements led to a *validation accuracy* of around 60%, while the *test accuracy* stands at a 57% on average.

Training progress and confusion matrix, along with a comparison with the baseline, will be discussed in Section 5.6.

5.5 Parameters and layout enhancements that did not provide benefit

Students chose the simplest network among many others that did not bring much of an improvement. The choice of the network has been determined by considering the accuracy, the complexity of the network, and the similarity with the baseline. If the network accuracy didn't improve, then

it shouldn't be more complex with added pieces. Among all those objectives, the student preferred networks with higher accuracy and complexity as lower as possible. In those cases where two networks were comparable, students preferred the network that was the most similar to the baseline with respect to learning parameters and layout.

With these fundamental goals in mind, the students discarded other configurations which did not result in significant improvements performance-wise, or introduced higher complexity, with less parameters shared between the improved network and the baseline.

In particular,

- increasing or lowering the size of *MiniBatch* has not been sufficient to improve performance – however, it changed significantly the behavior of the training phase, which required setting the *initial learning rate* to a different value;
- switching the optimization algorithm to *adam* has not been sufficient to show significant improvements, even after fine-tuning parameters. Students tried many different variants of the network with the adam optimization algorithm, however the results were not good enough to justify the switch;
- adding any convolutional layer did not result in performance improvements.

The students didn't investigate whether these improvements could bring performance improvements for any network – instead, the group only evaluated accuracy for the network after having applied all the changes discussed in the previous sections. Basically, none of the above modifications led to a significant performance improvement worth keeping.

5.6 Comparison with the baseline

The students collected pictures of training progress and relative confusion matrix of both *baseline network* (parameters in Section 4.9) and *improved network* (parameters in Section 5.4.3).

Figure 5 shows that classes that had been better predicted by the baseline network are still among those well predicted by the improved network too (for instance, *Street* and *TallBuilding*), while some other classes that were initially difficult have seen a huge improvement (*Mountain*, *LivingRoom*, and notably *Forest* class improved from a 5.7% to a 91.2%). Still, some classes look very hard to predict, for instance *Bedroom* class and *Industrial* ones.

In general, the improved network shows significant better performance when compared to the baseline network (an almost doubled accuracy on both validation and test set), at the cost of added complexity and some more training time required. Other variants of the improved network were tried by varying parameters and layout as described in Section 5.5, but no significant accuracy improvements justified the applied changes.

The overall code of the improved network can be read in `improved_network.m` file.

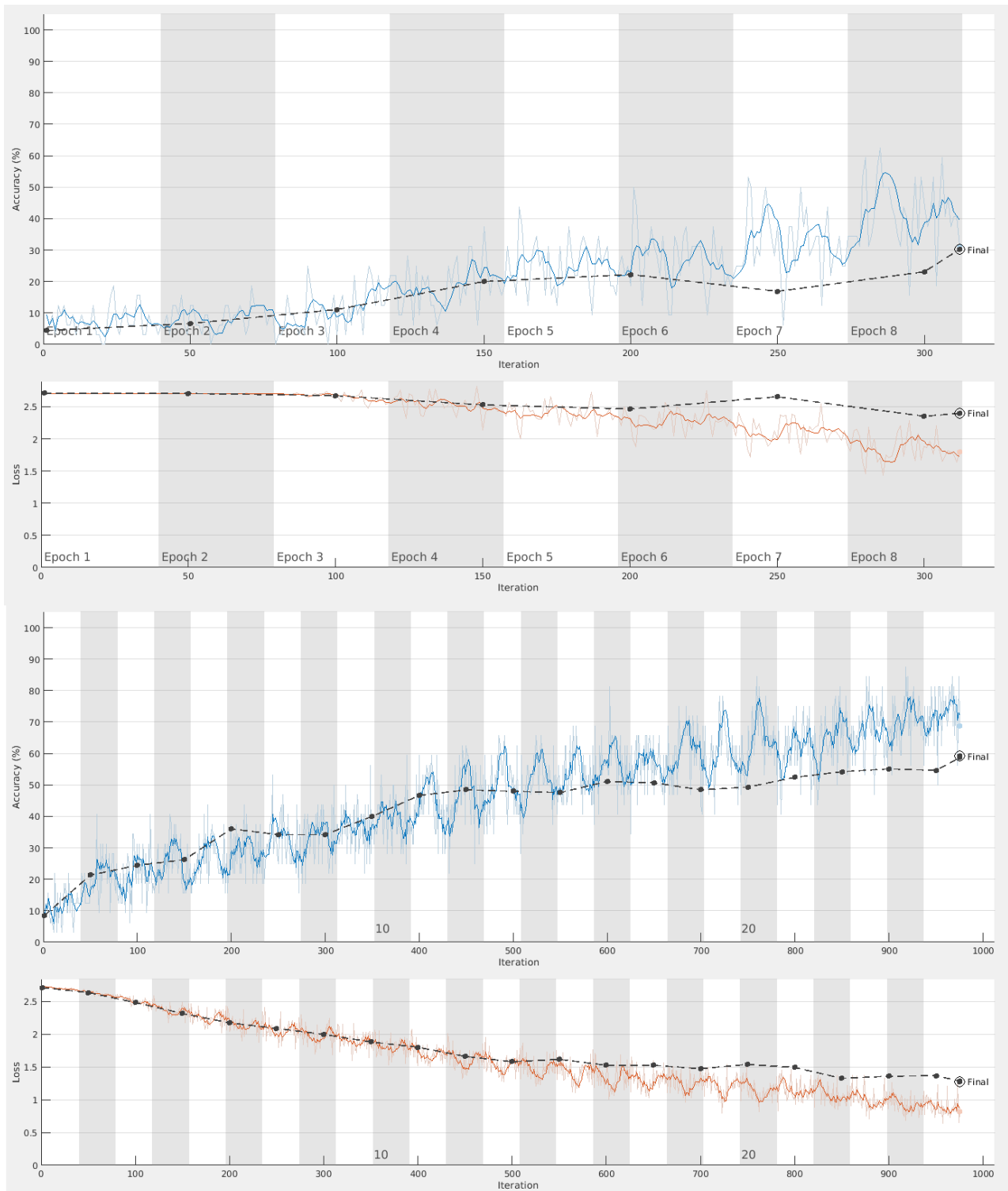


Figure 4: Comparison between training progress of the baseline network (on top) and training progress of the improved network (on bottom).

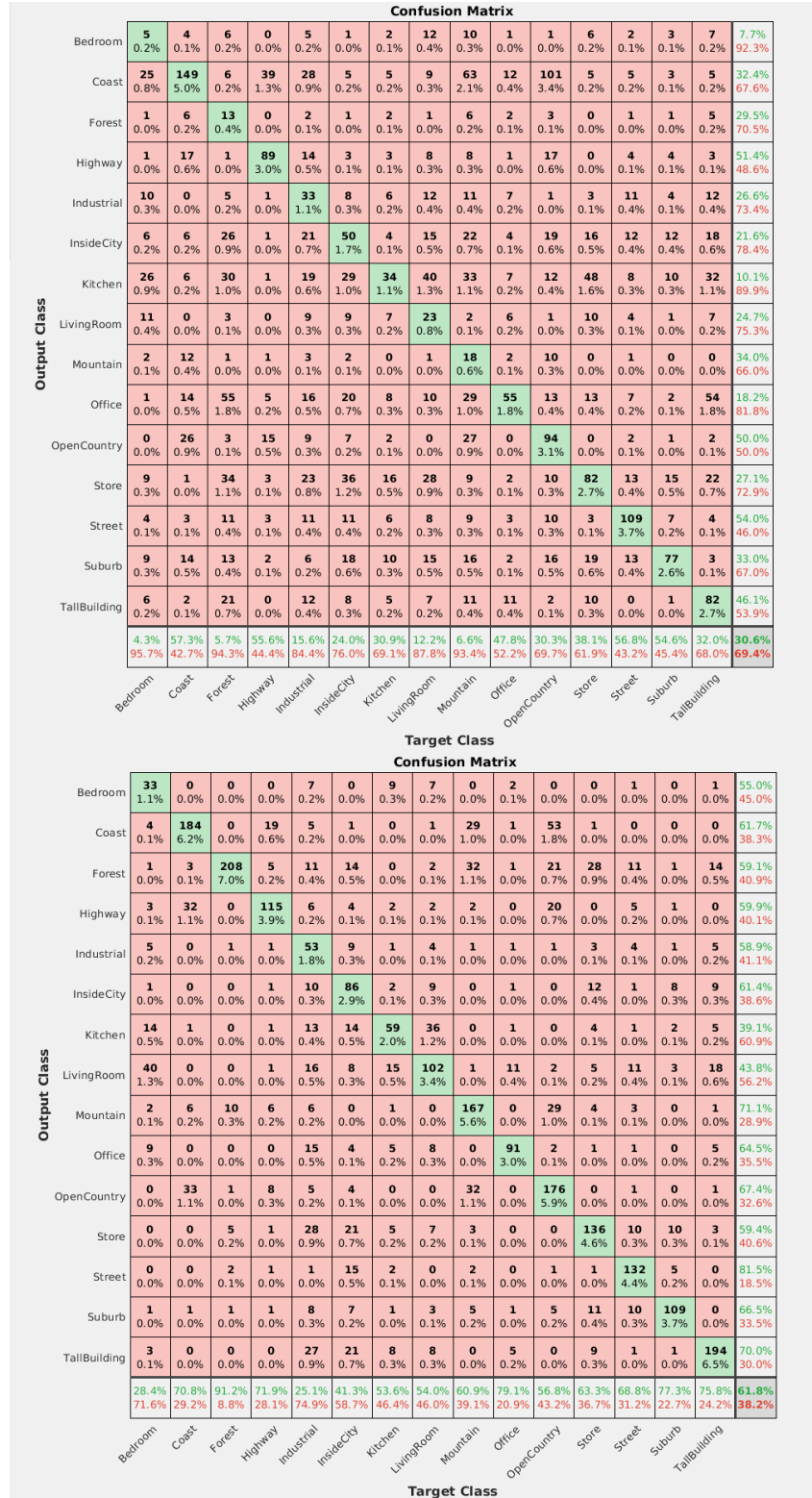


Figure 5: Comparison between confusion matrices of the baseline network (top) and the improved network (bottom).

5.7 Ensemble network based on previous improvements

The students implemented an **ensemble network**, composed of 5 networks of the same layout as in Section 5.4.3. In order to assign the class, a *voting* system has been built: each network of the ensemble casts a vote (a class), and the overall output is the *most common predicted class* among the ensemble.

The code related to the ensemble prediction is the following one,

```
% Evaluating ensemble network accuracy

EnsembleNum=5;
for i=1:EnsembleNum
    YPred = classify(net(i),imdsTest);
    for j=1:size(YPred)
        YPredicted(j,i) = YPred(j);
    end
end

categ = categorical(categories(YTest));
for i=1:size(YPred)
    [v, argmax] = max(countcats(YPredicted(i,:)));
    predicted_output(i) = categ(argmax);
end

accuracy = sum(predicted_output == YTest)/numel(YTest);
figure
plotconfusion(YTest, predicted_output')
```

where `net(i)` is a network belonging to an array of 5 trained networks. As a result, the above code will plot a confusion matrix for the ensemble network.

5.8 Justifying the added complexity

Ensemble training is very heavy on resources. Time training multiple networks is significantly higher than the time required for training a single network, let alone the time necessary for both training and classifying with an ensemble of networks. Hence, the students questioned themselves whether an ensemble of more networks (for instance, 10 networks or more) could bring any improvement with respect to an ensemble of 5 or less.

In order to answer the question, the students collected test accuracy of ensemble networks of different number B of networks, and compared them with the average of every single network in the ensemble. Each ensemble network has been trained separately, and results has been collected in the following Table 5,

Basically, training an ensemble of $B = 5$ networks should be enough to provide tangible benefits, with even higher improvements over the single networks when training a $B = 10$ ensemble network (although the improvements are not so pronounced).

B	Single networks average accuracy	Ensemble test set accuracy	Net improvement
3	0.5763	0.5950	1.87%
5	0.5863	0.6280	4.17%
7	0.5868	0.6340	4.72%
10	0.5841	0.6410	5.69%
15	0.5871	0.6470	5.99%

Table 5: Ensemble network accuracy evaluation.

Training the network with $B = 15$ only introduced a moderate improvement, thus such a high number of networks is not strictly required and would instead offer more disadvantages than advantages (for instance, training and classification time) with respect to an ensemble network having only 5 or 10 networks.

At every run, ensemble networks had shown more accuracy on the test set with respect to the single network. In addition, the single network performance varies a lot, largely depending on the randomness of the learning phase, with test set accuracy values ranging from 55% to over 60% – on contrary, the ensemble network accuracy has proven to be much more reliable over different runs. On this particular task, ensemble learning has proven to be effective – however, its increased complexity and time required for training and classification cannot be underestimated. A good trade-off between net improvement and increased complexity could be a choice of $B = 5$ networks for the ensemble network, which should bring the accuracy of the student’s proposed model well-above 60%.

The overall code of the ensemble network can be read in `ensemble_improved.m` file.

6 Adopting transfer learning with AlexNet

6.1 Transfer learning by freezing weights

In the third section, the students were required to adopt transfer learning based on a pre-trained network. The network of choice has been *AlexNet*: the first 22 layers of the network have been preserved, while only the last 3 were replaced by the last 3 layers of the improved network. Moreover, since the goal of transfer learning should be to adopt a pre-trained network with no further learning of its layers, the weights of the AlexNet layers have been frozen, preventing them to be modified during the learning phase.

The layout of the transfer learning network is described in Table 6:

#	Layer type	Size and parameters
1 : 22	AlexNet layers	weights freezed
23	Fully Connected	15
24	Softmax	softmax
25	Classification	crossentropyex

Table 6: Transfer Learning Layout

The code that produced the specified layout is the following one,

```

net = alexnet;

layersAlex = net.Layers(1:end-3);

layers = [
    layersAlex

    fullyConnectedLayer(15,'Name','fc_2',...
        'WeightsInitializer', @(sz) randn(sz)*0.01,...
        'BiasInitializer', @(sz) zeros(sz))

    softmaxLayer('Name','softmax')
    classificationLayer('Name','output')
];

layers(1:end-3) = freezeWeights(layers(1:end-3));

```

where the function `freezeWeights` has been obtained as-is from an online git repository (URL), and will be reported here:

```

function layers = freezeWeights(layers)
% layers = freezeWeights(layers) sets the learning rates of all the
% parameters of the layers in the layer array |layers| to zero.

for ii = 1:size(layers,1)
    props = properties(layers(ii));
    for p = 1:numel(props)
        propName = props{p};
        if ~isempty(regexpi(propName, 'LearnRateFactor$', 'once'))
            layers(ii).(propName) = 0;
        end
    end
end

end

```

6.1.1 Input layer

The students adopted the same data augmentation techniques that has been shown in Section 5.1. In order to preserve the input layer of the AlexNet (along with its subsequent structure), an input size of $227 \times 227 \times 3$ should be guaranteed. To get the correct input size students had to first modify the picture size accordingly, and then to implement a greyscale to RGB conversion. The group adopted the simple approach to replicate the greyscale image 3 times, creating the 3 color channels by simply providing the greyscale information to each channel.

The code relative to this part is the following one,

```

imds = imageDatastore(TrainDatasetPath, ...
    'IncludeSubfolders',true,'LabelSource','foldernames');
imds.ReadFcn = @(x)imresize(cat(3, imread(x), imread(x), imread(x)), [227 227]);
trainQuota=0.85;
[imdsTrain,imdsValidation] = splitEachLabel(imds,trainQuota,'randomize');
aug = imageDataAugmenter("RandXReflection",true);
imageSize = [227 227 3];
auimds = augmentedImageDatastore(imageSize,imdsTrain,'DataAugmentation',aug);

```

Notably, the main difference with respect to the previous sections is in the ReadFcn definition, where the image is first replicated into the 3 RGB channels, then rescaled to the desired size.

6.1.2 Training options

Different training options were explored. The group found out that good training options could be the following ones,

```

options = trainingOptions('sgdm', ...
    'MiniBatchSize',10, ...
    'MaxEpochs',6, ...
    'InitialLearnRate',1e-4, ...
    'ValidationData',imdsValidation, ...
    'ValidationFrequency',3, ...
    'Verbose',false, ...
    'Plots','training-progress');

```

6.1.3 Results

The training of the network provided an accuracy of 84.44% in the validation set, and of 86.4% in the test set, bringing overall accuracy well-above the 60% of the improved network illustrated in Section 5. Corresponding training progress and confusion matrix are provided in Figure 6 and Figure 7.

The code can be read in file `transfer_learning.m`.

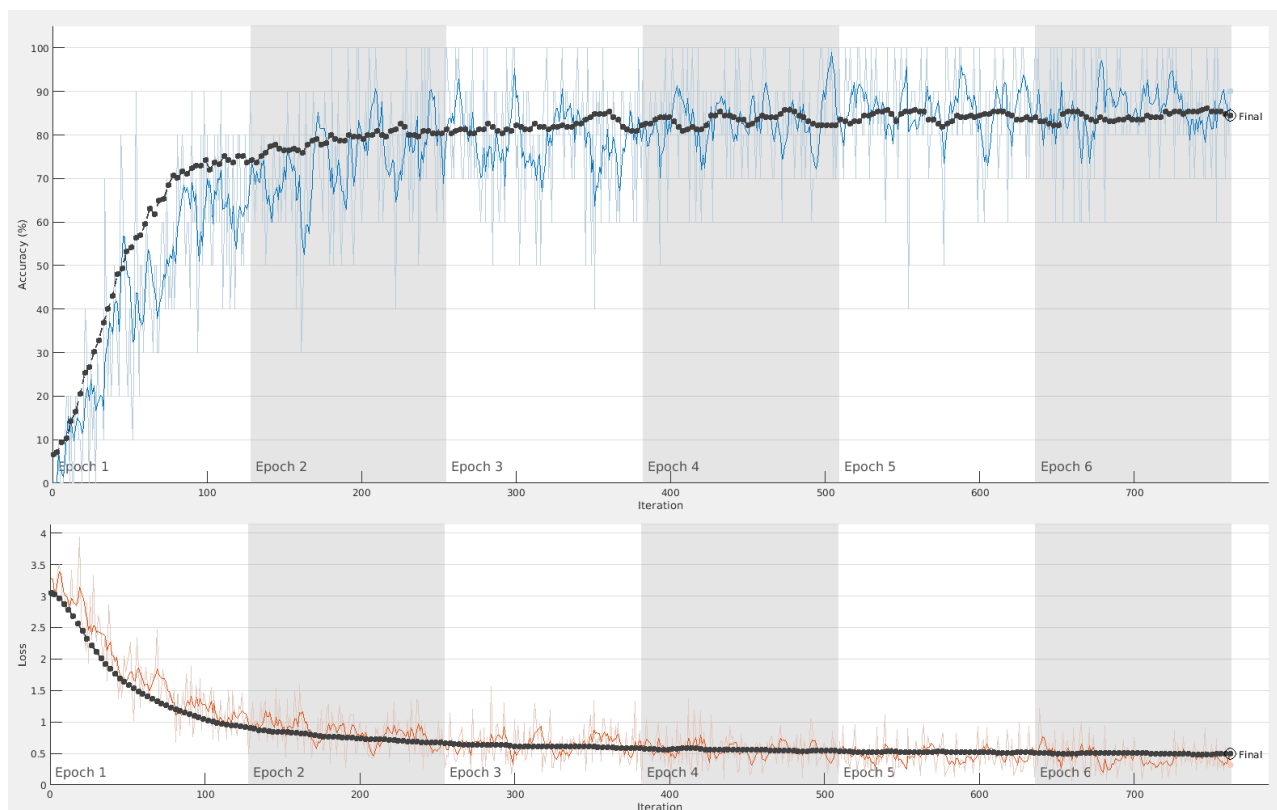


Figure 6: Progress of the transfer learning network, employing AlexNet with the last 3 layers modified.

Output Class		Confusion Matrix																	
		Bedroom	87 2.9%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	1 0.0%	3 0.1%	14 0.5%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	0 0.0%	0 0.0%	0 0.0%	80.6% 19.4%	
		Coast	0 0.0%	238 8.0%	0 0.0%	3 0.1%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	0 0.0%	39 1.3%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	84.1% 15.9%	
		Forest	0 0.0%	1 0.0%	221 7.4%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	4 0.1%	0 0.0%	17 0.6%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	90.6% 9.4%	
		Highway	0 0.0%	3 0.1%	0 0.0%	142 4.8%	3 0.1%	1 0.0%	0 0.0%	0 0.0%	2 0.1%	0 0.0%	4 0.1%	0 0.0%	6 0.2%	0 0.0%	0 0.0%	88.2% 11.8%	
		Industrial	0 0.0%	0 0.0%	0 0.0%	3 0.1%	161 5.4%	5 0.2%	0 0.0%	2 0.1%	0 0.0%	0 0.0%	2 0.1%	5 0.2%	4 0.1%	0 0.0%	18 0.6%	80.5% 19.5%	
		InsideCity	0 0.0%	0 0.0%	0 0.0%	2 0.1%	7 0.2%	177 5.9%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	1 0.0%	6 0.2%	5 0.2%	1 0.0%	5 0.2%	86.3% 13.7%	
		Kitchen	3 0.1%	0 0.0%	0 0.0%	0 0.0%	3 0.1%	2 0.1%	78 2.6%	6 0.2%	0 0.0%	2 0.1%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	82.1% 17.9%	
		LivingRoom	22 0.7%	0 0.0%	0 0.0%	2 0.1%	6 0.2%	0 0.0%	16 0.5%	151 5.1%	0 0.0%	4 0.1%	0 0.0%	5 0.2%	2 0.1%	2 0.1%	1 0.0%	71.6% 28.4%	
		Mountain	0 0.0%	1 0.0%	6 0.2%	0 0.0%	2 0.1%	0 0.0%	0 0.0%	0 0.0%	261 8.7%	0 0.0%	16 0.5%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	90.9% 9.1%	
		Office	4 0.1%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	1 0.0%	11 0.4%	11 0.4%	0 0.0%	108 3.6%	0 0.0%	5 0.2%	0 0.0%	0 0.0%	0 0.0%	76.6% 23.4%	
		OpenCountry	0 0.0%	17 0.6%	1 0.0%	2 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	5 0.2%	0 0.0%	231 7.7%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	90.2% 9.8%	
		Store	0 0.0%	0 0.0%	0 0.0%	1 0.0%	12 0.4%	3 0.1%	1 0.0%	3 0.1%	0 0.0%	1 0.0%	0 0.0%	185 6.2%	2 0.1%	0 0.0%	3 0.1%	87.7% 12.3%	
		Street	0 0.0%	0 0.0%	0 0.0%	3 0.1%	2 0.1%	9 0.3%	1 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	3 0.1%	172 5.8%	0 0.0%	1 0.0%	89.6% 10.4%	
		Suburb	0 0.0%	0 0.0%	0 0.0%	0 0.0%	4 0.1%	7 0.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	138 4.6%	0 0.0%	92.6% 7.4%	
		TallBuilding	0 0.0%	0 0.0%	0 0.0%	1 0.0%	9 0.3%	2 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	1 0.0%	0 0.0%	227 7.6%	93.8% 6.2%	
				75.0% 25.0%	91.5% 8.5%	96.9% 3.1%	88.8% 11.3%	76.3% 23.7%	85.1% 14.9%	70.9% 29.1%	79.9% 20.1%	95.3% 4.7%	93.9% 6.1%	74.5% 25.5%	86.0% 14.0%	89.6% 10.4%	97.9% 2.1%	88.7% 11.3%	86.3% 13.7%
				Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding	
				Target Class															

Figure 7: Confusion matrix of the transfer learning network with last layers modified.

6.2 AlexNet as feature extractor

AlexNet could be employed as a *feature extractor*. By accessing the activation of an intermediate layer, it should be possible to extract the features that could be used to train a SVM classifier. The students extracted the activation of the latest convolutional layer in AlexNet ('fc7') and trained a multiclass SVM classifier employing Error-Correcting Output Codes.

The code should be the following,

```
net = alexnet;

featuresTrain = activations(net, imdsTrain, 'fc7', 'OutputAs', 'rows');
YTrain = imdsTrain.Labels;
classifier = fitcecoc(featuresTrain, YTrain);

TestDatasetPath = fullfile('dataset', 'test');
imdsTest = imageDatastore(TestDatasetPath, ...
    'IncludeSubfolders', true, 'LabelSource', 'foldernames');
imdsTest.ReadFcn = @(x) imresize(cat(3, imread(x), imread(x), imread(x)), [227 227]);
YTest = imdsTest.Labels;
featuresTest = activations(net, imdsTest, 'fc7', 'OutputAs', 'rows');
YPredicted = predict(classifier, featuresTest);
figure
plotconfusion(YTest, YPredicted)
```

where the last instruction should plot a confusion matrix such as the one in Figure 8.

The test set accuracy stands at a 87.6%, a result that is comparable to the one obtained in Section 6.1.

The code can be read in file `transfer_learning_svm.m`.

		Confusion Matrix															
Output Class	Bedroom	96 3.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	4 0.1%	23 0.8%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	76.8% 23.2%	
	Coast	0 0.0%	224 7.5%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	3 0.1%	0 0.0%	16 0.5%	0 0.0%	0 0.0%	0 0.0%	91.8% 8.2%	
	Forest	0 0.0%	1 0.0%	218 7.3%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	0 0.0%	5 0.2%	0 0.0%	0 0.0%	1 0.0%	95.6% 4.4%	
	Highway	0 0.0%	3 0.1%	0 0.0%	145 4.9%	2 0.1%	1 0.0%	0 0.0%	0 0.0%	4 0.1%	0 0.0%	2 0.1%	0 0.0%	6 0.2%	0 0.0%	89.0% 11.0%	
	Industrial	0 0.0%	0 0.0%	0 0.0%	2 0.1%	171 5.7%	11 0.4%	0 0.0%	4 0.1%	1 0.0%	0 0.0%	1 0.0%	9 0.3%	5 0.2%	1 0.0%	76.0% 24.0%	
	InsideCity	0 0.0%	0 0.0%	0 0.0%	1 0.0%	9 0.3%	178 6.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	6 0.2%	2 0.1%	2 0.1%	88.1% 11.9%	
	Kitchen	3 0.1%	0 0.0%	0 0.0%	0 0.0%	5 0.2%	4 0.1%	88 2.9%	13 0.4%	0 0.0%	2 0.1%	0 0.0%	2 0.1%	1 0.0%	0 0.0%	74.6% 25.4%	
	LivingRoom	17 0.6%	0 0.0%	0 0.0%	2 0.1%	3 0.1%	2 0.1%	12 0.4%	141 4.7%	0 0.0%	9 0.3%	0 0.0%	4 0.1%	2 0.1%	2 0.1%	72.7% 27.3%	
	Mountain	0 0.0%	2 0.1%	1 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	244 8.2%	0 0.0%	5 0.2%	0 0.0%	0 0.0%	0 0.0%	96.4% 3.6%	
	Office	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	5 0.2%	3 0.1%	0 0.0%	104 3.5%	0 0.0%	4 0.1%	0 0.0%	0 0.0%	88.9% 11.1%	
	OpenCountry	0 0.0%	30 1.0%	9 0.3%	6 0.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	20 0.7%	0 0.0%	281 9.4%	0 0.0%	0 0.0%	1 0.0%	81.0% 19.0%	
	Store	0 0.0%	0 0.0%	0 0.0%	2 0.1%	9 0.3%	1 0.0%	1 0.0%	4 0.1%	0 0.0%	0 0.0%	0 0.0%	187 6.3%	1 0.0%	0 0.0%	90.8% 9.2%	
	Street	0 0.0%	0 0.0%	0 0.0%	1 0.0%	4 0.1%	5 0.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	174 5.8%	0 0.0%	94.1% 5.9%	
	Suburb	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	3 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	134 4.5%	96.4% 3.6%	
	TallBuilding	0 0.0%	0 0.0%	0 0.0%	0 0.0%	5 0.2%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	1 0.0%	0 0.0%	231 7.7%	96.7% 3.3%
		82.8% 17.2%	86.2% 13.8%	95.6% 4.4%	90.6% 9.4%	81.0% 19.0%	85.6% 14.4%	80.0% 20.0%	74.6% 25.4%	89.1% 10.9%	90.4% 9.6%	90.6% 9.4%	87.0% 13.0%	90.6% 9.4%	95.0% 5.0%	90.2% 9.8%	87.6% 12.4%
		Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding	
		Target Class															

Figure 8: Confusion matrix of the transfer learning network. This time, AlexNet worked as a feature extractor (up to the last convolutional layer), while the classification task is reserved to a SVM that has been trained with the extracted features.

7 Summary of results

The students implemented a total of 5 notable Convolutional Neural Networks,

1. the **baseline network** in Section 4;
2. the **improved network** in Section 5;
3. the **ensemble network** in Section 5.7;
4. a network implementing **transfer learning** by substituting the last 3 layers of AlexNet (Section 6.1);
5. a network implementing **transfer learning** by extracting activations of the last convolutional layer of AlexNet and training a SVM (Section 6.2).

The baseline network mainly served as a comparison and a starting point. The improved network almost doubled the accuracy of the baseline network, by only implementing some crucial changes (notably, data augmentation and addition of batch normalization layers were enough to improve drastically the performance of the network. Subsequent improvements refined the learning phase and granted enough benefits to double the performance of the baseline). Then, the improved network was adopted to build an ensemble network composed of 5 and more of those networks. Each network was asked to classify an image; the output of the ensemble network was the most commonly predicted class. The approach granted an accuracy improvement of 4.17% and much more stability of the results in the case of the improved network. The best network out of the 3 first implemented by the students has been the ensemble network – however, the training time has been significantly increased by the necessity of learning multiple networks one after another.

For the last 2 networks of the list, the group started from a pre-trained network (AlexNet) and only implemented the last portion of it, along with the testing apparatus. The first approach was by substituting the last layers with the last ones from the improved network; the second approach, instead, extracted features from the last convolutional layer and trained a SVM for classification. By adopting a well-established network, both approaches ended up being successful, with a test set accuracy of, respectively, 86.4% and 87.6%. Those two networks fully outclassed the previous ones, showing that the transfer learning approach is not only feasible, but could be a valid and quick technique to obtain a high-performance Convolutional Neural Network.

References

- [1] [Lazebnik et al., 2006] Lazebnik, S., Schmid, C., and Ponce, J. (2006). Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), volume 2, pages 2169–2178.