

TDT4265: Computer Vision and Deep Learning

Assignment 4

Håkon Hukkelås
hakon.hukkelas@ntnu.no
Department of Computer Science
Norwegian University of Science and Technology (NTNU)

March 12, 2020

- **Delivery deadline: Monday, March 23, 2020, by 23:59.**
- **This assignment count towards 6% of your final grade.**
- You can work on your own or in groups of up to 2 people.
- Upload your code as a single ZIP file.
- Upload your report as a single PDF file to blackboard.
- You are required to use python3 to finish the programming assignments.
- The delivered code is taken into account with the evaluation. Ensure your code is well documented and as readable as possible.

Introduction

In this assignment, we will take a look at object detection with convolutional neural networks. First, we will go through a common metric used to evaluate object detection models: mean average precision. Furthermore, we will take a deep dive into the SSD300 architecture, a lightweight real-time object detection model.

With this assignment, we provide you starter code for the programming tasks. You can download this from:

<https://github.com/hukkelas/TDT4265-StarterCode>.

To set up your environment, follow the guide in the Github repo:

https://github.com/hukkelas/TDT4265-StarterCode/blob/master/python_setup_instructions.md

Recommended Readings

We recommend you to review the lectures slides to get a brief overview of convolutional neural networks before starting this assignment. In addition, these are some recommended readings to get you started on the assignment.

1. [Section 4.2 of The PASCAL Visual Object Classes \(VOC\) Challenge](#). This is the official description of how to calculate the mean average precision for object detection tasks.
2. [Jonathan Hui's blog post about mean average precision](#). This blog post is very short, but gives a very simple explanation of mean average precision.
3. [SSD Blog Post by Jonathan Hui](#): We recommend everyone to read this before starting on the SSD tasks.
4. [SSD: Single Shot MultiBox Detector](#). We recommend you to read the blog post before reading the paper.

Delivery

We ask you to follow these guidelines:

- **Report:** Deliver your answers as a **single PDF file**. Include all tasks in the report, and mark it clearly with the task you are answering (Task 1.1, Task1.2, Task 2.1a etc).
- **Plots in report:** For the plots in the report, ensure that they are large and easily readable. You might want to use the "ylim" function in the matplotlib package to "zoom" in on your plots. Label the different graphs such that it is easy for us to see which graphs correspond to the train, validation and test set.
- **Source code:** Upload your code as a zip file. In the assignment starter code, we have included a script (`create_submission_zip.py`) to create your delivery zip. **Please use this**, as this will structure the zipfile as we expect. (Run this from the same folder as all the python files).
To use the script, simply run: `python3 create_submission_zip.py`
- **Upload to blackboard:** Upload the ZIP file with your source code and the report to blackboard before the delivery deadline.

Late Delivery & Other Penalties

We have a strict policy on late deliveries. Any delivery submitted after the deadline will be limited in maximum score which you can achieve on the assignment. See "Course Information" on blackboard for more information.

Also, if you fail to follow our delivery guidelines we will **subtract 0.2 points** from your final score. This includes if you do not deliver your report as a PDF file, if you don't deliver your code as .ZIP file, or if you don't use the `create_submission_zip.py` to generate your .zip file.

1 Object Detection Metrics (0.5 point)

For all previous assignments, we have done classification of images. When measuring the performance of image classification architectures we often use accuracy. For object detection, measuring performance is a more complicated task. First, we need to ensure that we detect the object correctly with a good bounding box. Then, when we have detected the object, we want to make sure that the detected object is classified correctly. When the PASCAL VOC dataset was proposed, they introduced a standard metric to measure the performance of object detection architectures. This metric was the *mean average precision* (*mAP*). We recommend you to take a look at the recommended resources about mean average precision.

In your report, please:

- (a) [0.1pts] Explain what the Intersection over Union is and how we can find it for two bounding boxes. Illustrate it with a drawing.
- (b) [0.1pts] Write down the equation of precision and recall, and shortly explain what a true positive and false positive is.
- (c) [0.3pts] Given the following precision and recall curve for the two classes, what is the mean average precision?

Precision and recall curve for class 1:

Precision₁ = [1.0, 1.0, 1.0, 0.5, 0.20]

Recall₁ = [0.05, 0.1, 0.4, 0.7, 1.0]

Precision and recall curve for class 2:

Precision₂ = [1.0, 0.80, 0.60, 0.5, 0.20]

Recall₂ = [0.3, 0.4, 0.5, 0.7, 1.0]

Hint: To calculate this, find the precision for the following recall levels: 0.0, 0.1, 0.2, ... 0.9, 1.0.

2 Implementing Mean Average Precision (2 points)

Now, we will implement mean average precision for a single class. In the assignment files, there is a .json file for both a set of predicted bounding boxes and a set of ground truth bounding boxes. With the assignment files, you will find a file called task2.py and this is where all the code should go for this task. Also, we have provided you a set of tests to simplify the debugging steps for this task. For each subtask you will finish, you can run tests.py and it will check if the function is correct. These tests are written with simple checks and we can't guarantee that they cover every edge case. If you want to implement this in another language you are free to do so, but we will not provide tests for other languages.

Note, for this task we will implement average precision for a single class. If you want to extend this code for a multi-class object detection algorithm, you find the average precision for each class and take the mean over all of these to get the mean average precision.

Implement the following:

- (a) [0.3pts] Implement a function that takes two bounding boxes, then finds the intersection over union. This is the function `calculate_iou` in task2.py
- (b) [0.2pts] Implement a function that computes the precision given the number of true positives, false positives, and false negatives. This is the function `calculate_precision` in task2.py.
Implement a function that computes the recall given the number of true positives, false positives, and false negatives. This is the function `calculate_recall` in task2.py
- (c) [0.5pts] Implement a function that takes in a set of predicted bounding boxes, a set of ground truth bounding boxes, and a given intersection of union threshold, then finds the optimal match for each bounding box. This is the function `get_all_box_matches` in task2.py.
Note, each box can only have a single match and it should be matched with the bounding box with the highest intersection of union ratio. If there is no match that has a higher IoU than the IoU threshold then the bounding box has no match.
- (d) [0.4pts] Implement a function that calculates the true positives, false positives, and false negatives for a single image. Furthermore, implement a function that finds the precision and recall for all images in a dataset. Finally, implement a function that computes the precision recall curve. This are the functions `calculate_individual_image_result`, `calculate_precision_recall_all_images`, and `get_precision_recall_curve` in task2.py.
- (e) [0.5pts] Implement a function that computes the mean average precision given a list of precisions and recalls. This is the function `calculate_mean_average_precision` in task2.py
Note, do not change the recall levels used to compute the mAP. This is what is used in standard object detection benchmarks.
If you run task2.py as a python file, the printed mean average precision over the dataset should be 0.9066 if you use a IoU threshold of 0.5.
- (f) [0.1pts] **In report:** Put the plot of your final precision-recall curve in your report.
If you run task2.py, it will save the precision-recall curve to precision_recall_curve.png

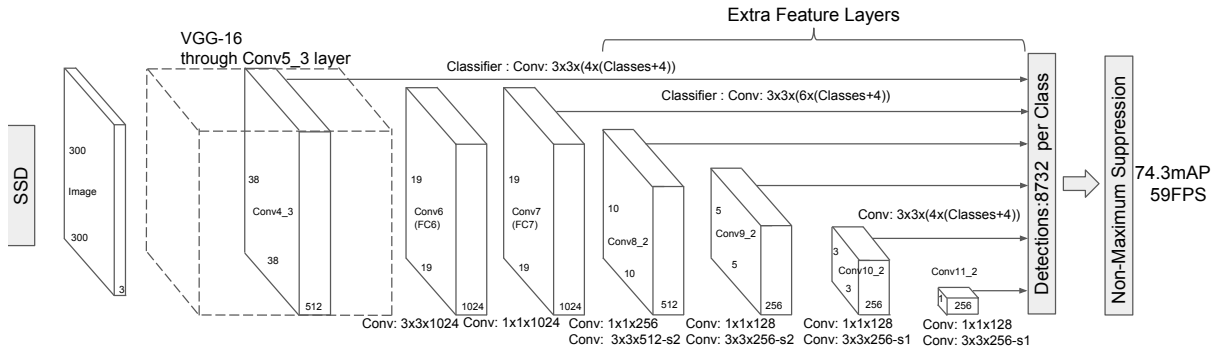


Figure 1: The SSD architecture presented in the original paper.

3 Theory [0.5 point]

All of these questions can be answered by reading the SSD paper.

1. [0.1pts] The SSD architecture produces a fixed-size number of bounding boxes and a score for each bounding box. After this, the final step is to filter out overlapping boxes, what is this called?
2. [0.1pts] The SSD architecture predicts bounding boxes at multiple scales to enable the network to detect objects of different sizes.

Is the following true or false: Predictions from the deeper layers in SSD are responsible to detect small objects.

3. [0.2pts] SSD use k number of bounding boxes with different aspect ratios at each spatial location in a feature map to predict c class scores and 4 offsets relative to the original box shape.

Why do they use different bounding box aspect ratios at the same spatial location?

4. [0.1pts] What is the main difference between SSD and YOLOv1/v2 (The YOLO version they refer to in the SSD paper)?

Table 1: Illustration of a custom SSD CNN backbone. The number of filters specifies the number of filters/kernels in a convolutional layer. Each convolutional layer has a filter size of 3×3 and padding of 1. Each MaxPool2D layer has a stride of 2 and a kernel size of 2×2 . Note that `output_channels[i]` is a variable that you can set in the config file (By default this is: [128, 256, 128, 128, 64, 64]). * **The last convolution** should have a stride of 1, padding of 0, and a kernel size of 3×3 .

Is Output	Layer Type	Number of Filters	Stride
Yes - Resolution: 38×38	Conv2D	32	1
	MaxPool2D	–	2
	ReLU	–	–
	Conv2D	64	1
	MaxPool2D	–	2
	ReLU	–	–
	Conv2D	64	1
	ReLU	–	–
	Conv2D	<code>output_channels[0]</code>	2
Yes - Resolution: 19×19	ReLU	–	–
	Conv2D	128	1
	ReLU	–	–
	Conv2D	<code>output_channels[1]</code>	2
Yes - Resolution: 9×9	ReLU	–	–
	Conv2D	256	1
	ReLU	–	–
	Conv2D	<code>output_channels[2]</code>	2
Yes - Resolution: 5×5	ReLU	–	–
	Conv2D	128	1
	ReLU	–	–
	Conv2D	<code>output_channels[3]</code>	2
Yes - Resolution: 3×3	ReLU	–	–
	Conv2D	128	1
	ReLU	–	–
	Conv2D	<code>output_channels[4]</code>	2
Yes - Resolution: 1×1 *	ReLU	–	–
	Conv2D	128	1
	ReLU	–	–
	Conv2D	<code>output_channels[5]</code>	1

4 Single Shot Detector [3 points]

For this task, we provide a large code base with an SSD implementation. Most of the code is finished, except the backbone network. Your task is to implement the backbone which should output the **six** different feature banks that SSD use to predict bounding boxes (see [Figure 1](#)). We will keep the resolution of each feature bank; however, we will reduce the number of feature maps per bank. ¹ [Table 1](#) shows the backbone network you will implement.

This code base is significantly larger than what we've given you before. To get you started, take a look in the [readme.md](#) file in the repository:

Note that for this task we will only use a train and validation set (no test set).

- (a) [1pts] Implement the CNN in [Table 1](#). The only code required to change is located in `ssd/modeling/backbone/basic.py`.

¹Note that a feature bank consists of several feature maps. A single feature map is the output from a single convolution filter. A feature bank is the combined feature map of several filters in a single layer.

- (b) [0.5pts] Train your model on the MNIST Object detection dataset. For information about the dataset, see: <https://github.com/hukkelas/MNIST-ObjectDetection>.

You should expect a mean average precision (mAP) of about 75-77% if you train for 6000 iterations. In your report, include a plot of "total_loss" (You can take a screenshot from the tensorboard, or use the provided jupyter notebook).

Fianlly, report your mean average precision (mAP).

- (c) [1pts] Use what you learned from assignment 3 to improve your model. Reach a mean average precision of 85% within 10K gradient descent iterations (or get as close as possible - we will give partial points).

In your report, describe what improvements you did to reach the accuracy and report the final mAP.

Some tips for improving your model:

- Look back at assignment 3 and think about what improved your model. These additions will probably improve this model as well!
- You can change other parts of the code as well (for example, augmentation, optimizers, learning rate, etc.)
- Be careful to not add too much data augmentation. Our reference implementation used no data augmentation. (Some might help though!)
- The model in Table 1 has somewhere in the neighbourhood of 2M parameters. The original backbone (VGG) had 130M. Maybe adding more layers or filters per layer might work?

- (d) [0.2pts] **(HARD)** Use whatever means possible to reach 90% mAP on the validation set within 15K gradient descent iterations .

- (e) [0.1pts] Try out your model on a couple of images. Run demo.py on the images in the demo/mnist folder and include the classified images in your report.

Do this by typing the following:

```
python demo.py --images_dir demo/mnist --dataset_type mnist configs/mnist.yaml
```

Was there any digits your model was not able to detect?

- (f) [0.2pts] This last task is to prepare you for the project and test a model on the PASCAL VOC dataset. VOC consist of 21 different classes and we will train a SSD network with VGG16 as the backbone. This backbone is already pre-trained on the ImageNet dataset (classification).

We've included documentation on how to download the VOC dataset in the readme.md file in the starter code.

Use the config file "vgg_ssd300_voc0712.yaml" to train a VGG16 model on the VOC dataset. Train the model for 5000 gradient descent iterations (the same as "iter" that prints in the starter code).

Test your model on a couple of images, by running the demo.py file:

```
python demo.py --images_dir=demo/voc --dataset_type=voc configs/vgg_ssd300_voc0712.yaml
```

In you report, include the plot of "total_loss", your final mAP on the validation set and the images you tested with demo.py.

Note that no modification to the code should be required to finish this task.