

IN5400

Alex

Week 07: Fine Tuning of neural networks

[The following notes are compiled from various sources such as textbooks, lecture materials, Web resources and are shared for academic purposes only, intended for use by students registered for a specific course. In the interest of brevity, every source is not cited. The compiler of these notes gratefully acknowledges all such sources.]

1 two exemplary ways of fine-tuning, setup with a separate validation set

- take the 102 class flowers dataset <https://www.dropbox.com/s/dx5ikzvrlqs2bbw/102flowersn.zip?dl=0> and write a dataset class which can work with a train/val/test split for it. The difference to the 102 flowers from oxford is that I provided in it for you train/val/test split for your convenience (bowing politely in front of the students).
- take any deep network you like which has pretrained weights (a smaller resnet, a smaller densenet are training fast)
- train a deep neural network in three different modes
 - (A) once without loading weights and training all layers.
 - (B) once with loading model weights before training and training all layers,
 - (C) once with loading model weights before training and training only the last **one** trainable layer (note: for quite some problems, the approach B is better than C)

For each of these 2 modes select the best epoch by the performance of the model on the validation set. Typically less than 20 epochs should suffice for training when using finetuning. You can run also optionally a selection over a few learning rates, if you use a GPU.

What do you need to do for steps when you start with a code like the MNIST training code?

- write a new dataloader for your training dataset, the flowers
- adjust paths for data (and if necessary for label paths/files or files determining splits into train/val/test)
- write the code so that it works with paths relative to the directory of your script
- as you did not have data augmentation so far (next lecture), you can simply use the following:

```
data_transforms = {
    'train': transforms.Compose([
        transforms.Resize(256),
        transforms.RandomCrop(224),
        transforms.RandomHorizontalFlip(),
        transforms.ToTensor(),
        transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
    ]),
    'val': transforms.Compose([
        transforms.Resize(224),
        transforms.CenterCrop(224),
        transforms.RandomHorizontalFlip(),
        transforms.ToTensor(),
        transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
    ]),
}
```

I am aware that I am using here the imagenet mean and the imagenet standard deviation.

- use some deep learning model from the model zoo, load its weights before training for settings B and C
- adjust the number of classes in the last linear/dense layer
- in the optimizer provide the proper parameters for training for settings A,B,C
- collect and plot curves of the training loss, the validation set loss and the validation set accuracy. also print the final test accuracy of the selected model
- A note: Calling a model constructor with `pretrained=True` does not tell you what really goes on when one. Check <https://github.com/pytorch/vision/blob/master/torchvision/models/resnet.py> to see what routine is used to load a model.

1.1 Bonus Task

Try out the effect of mixup training on your problem:

<https://arxiv.org/abs/1710.09412>. The paper is good to read and it is not too complicated to implement that one. I have no guess how much it can bring here because the flowers are a relatively simple problem. It is good to know this kind of trick anyway.

2 Some theory

- You are given a 2-dimensional convolution with feature map input size $(78, 84)$. When using a kernel of size $(5, 5)$ and stride 3 with padding of 2, what will be the spatial size of the feature map which is the output of the convolution? Note that the *spatial* size does not depend on the number of input or output channels.
- You are given a 1-dimensional convolution. When using a kernel of size 9 and stride 3 with padding 1, which spatial input size do you need to have, so that you have a spatial output size of 16?
- You are given a 2-dimensional convolution. When using a kernel of size $(3, 5)$ and stride 2 with padding of 0, what will be the spatial size of the feature map which is the output of the convolution, which spatial input size do you need to have, so that you have a spatial output size of $(128, 96)$?

3 Some more theory

How many trainable parameters are in

- a 2-D convolutional layer with input $(32, 19, 19)$, kernel size $(7, 7)$, stride 3, 64 output channels?
- a 2-D convolutional layer with input $(512, 25, 25)$, kernel size $(1, 1)$, stride 1, 128 output channels?
- how many multiplications and how many additions are performed in the first case above?