

Mini-Project Report

Hao Wang, Puda Zhao, Yinhong Qin

New York University
Deep Learning Fall 2022 ECE-GY 7123
hw2671@nyu.edu, pz2078@nyu.edu, yq2021@nyu.edu

Abstract

Deep neural network model can estimate more complicated function theoretically but it will suffer from the gradient explosion and vanishing as well as the exponentially growth of the number of parameters. Residual neural network (ResNet) can improve these problems and help researchers build very deep network. In this mini-project, we propose a image classification model based on ResNet model with the total number of parameters less than 5M. Our model was trained and tested on CIFAR-10 dataset and can reach 93.8% accuracy. We also made several variations of changes to select the final model with best testing performance on this dataset.

Overview

We trained a ResNet architecture to classify images from the CIFAR10 dataset with maximum around 0.94 accuracy. We used several useful techniques in this project to improve our model performance and reduced the training time, i.e., Data normalization, Data augmentation, Learning rate scheduling, etc. We will explain in detail in the sub-sequent section.

We tried several different Block Number combinations and found that the closer the change of Block Number is to the input, the smaller the impact on the number of parameters, and the closer to the output, the greater the impact on the number of parameters.

We adjusted the parameters on the default ResNet-18 model, and achieved better accuracy than the default model with a more simplified model. The training time for 70 epochs is about 25 minutes.

Our codebase can be accessed through this GitHub repository: <https://github.com/KevinQyhNYU/Deep-learning-Mini-project/tree/master>

It is modified based some online examples: For code source part, we checked and did modifications based on: <https://www.kaggle.com/code/datajameson/cifar-10-object-recognition-resnet-acc-94>.

For ResNet design part, our model is based on: <https://github.com/kuangliu/pytorch-cifar/blob/master/models/resnet.py>.

Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

```
Data Transformation (Data Normalisation and Data Augmentation)

1 state= ((0.4914,0.4822,0.4465),(0.2021,0.1964,0.2019)) mean and std
2 train_tf= tt.Compose([tt.RandomCrop(32, padding=4, padding_mode='reflect'), # transformation of data together
3                       tt.RandomHorizontalFlip(),
4                       tt.ToTensor(),tt.Normalize(*state,inplace=True)])
5 valid_tf= tt.Compose([tt.ToTensor(),tt.Normalize(*state)])
✓ 6h

1 # Image transformation
2 train_ds = datasets.CIFAR10(root='./data', train=True, download=True, transform=train_tf) # Data augmentation is only done on training images
3 valid_ds = datasets.CIFAR10(root='./data', train=False, download=True, transform=valid_tf)
4
✓ 1h
Files already downloaded and verified
Files already downloaded and verified
```

Figure 1: Data normalization and augmentation

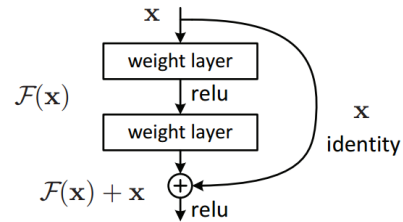


Figure 2: simple model of ResNet

Methodology

Data normalization: We normalized the image tensors by subtracting the mean and dividing by the standard deviation of pixels across each channel. Normalizing the data prevents the pixel values from any one channel from disproportionately affecting the losses and gradients.

Data augmentation: We applied random transformations while loading images from the training dataset. Specifically, we will pad each image by 4 pixels, and then take a random crop of size 32 x 32 pixels, and then flip the image horizontally with a 0.5 probability.

Residual connections: One of the key changes to our CNN model was the addition of the residual block, which adds the original input back to the output feature map obtained by passing the input through one or more convolutional layers. To seamlessly use a GPU, if one is available, we define a couple of helper functions.

Batch normalization: After each convolutional layer, we added a batch normalization layer, which normalizes the outputs of the previous layer. This is somewhat similar to data normalization, except it's applied to the outputs of a

GPU Helper Functions

To use a GPU, if one is available, we define a couple of helper functions (get_default_device & to_device) and a helper class DeviceDataLoader to move our

```

1 def get_default_device():
2     """Pick one if available, else cpu"""
3     if torch.cuda.is_available():
4         return torch.device('cuda')
5     else:
6         return torch.device('cpu')
7
8 def to_device(data, device):
9     """Move tensor(s) to chosen device"""
10    if isinstance(data, (list, tuple)):
11        return [to_device(x, device) for x in data]
12    return data.to(device, non_blocking=True)
13
14 class DeviceDataLoader():
15     """Wrap a dataloader to move data to a device"""
16     def __init__(self, dl, device):
17         self.dl = dl
18         self.device = device
19
20     def __iter__(self):
21         """Yield a batch of data after moving it to device"""
22         for b in self.dl:
23             yield to_device(b, self.device)
24
25     def __len__(self):
26         """Number of batches"""
27         return len(self.dl)
28
29 ✓ 66

```

Figure 3: GPU functions

layer, and the mean and standard deviation are learned parameters.

Learning rate scheduling: Instead of using a fixed learning rate, we will use a learning rate scheduler, which will change the learning rate after every batch of training. There are many strategies for varying the learning rate during training, and we used the "One Cycle Learning Rate Policy".

Weight Decay We added weight decay to the optimizer, yet another regularization technique which prevents the weights from becoming too large by adding an additional term to the loss function.

Gradient clipping: We also added gradient clipping, which helps limit the values of gradients to a small range to prevent undesirable changes in model parameters due to large gradient values during training.

Results

We designed a 4-layer ResNet with block numbers equal to [2, 1, 1, 1]. We tried different cases by setting different weight-decay and finally reached maximum 0.9381 accuracy with total 4977226 of parameters.

```

[26] epochs = 70
max_lr = 0.01
grad_clip = 0.1
weight_decay = 1e-6
opt_func = torch.optim.Adam

[26] %%time
history += fit_one_cycle(epochs, max_lr, model, train_dl, valid_dl,
                        grad_clip=grad_clip,
                        weight_decay=weight_decay,
                        opt_func=opt_func)

```

Figure 4: Settings

Figure 5 to 8 show the information about training process and testing results of this model. Based on the loss curve of the training and validation process, we can notice that this model does not have the overfitting issue on the CIFAR-10 dataset.

Figure 9 is a visualization result showing how the learning rate was changing over the training epochs. We can notice

Layer (type)	Output Shape	Param #
Conv2d-1	[-1, 64, 32, 32]	1,728
BatchNorm2d-2	[-1, 64, 32, 32]	128
Conv2d-3	[-1, 64, 32, 32]	36,864
BatchNorm2d-4	[-1, 64, 32, 32]	128
Conv2d-5	[-1, 64, 32, 32]	36,864
BatchNorm2d-6	[-1, 64, 32, 32]	128
BasicBlock-7	[-1, 64, 32, 32]	0
Conv2d-8	[-1, 64, 32, 32]	36,864
BatchNorm2d-9	[-1, 64, 32, 32]	128
Conv2d-10	[-1, 64, 32, 32]	36,864
BatchNorm2d-11	[-1, 64, 32, 32]	128
BasicBlock-12	[-1, 64, 32, 32]	0
Conv2d-13	[-1, 128, 16, 16]	73,728
BatchNorm2d-14	[-1, 128, 16, 16]	256
Conv2d-15	[-1, 128, 16, 16]	147,456
BatchNorm2d-16	[-1, 128, 16, 16]	256
Conv2d-17	[-1, 128, 16, 16]	8,192
BatchNorm2d-18	[-1, 128, 16, 16]	256
BasicBlock-19	[-1, 128, 16, 16]	0
Conv2d-20	[-1, 256, 8, 8]	294,512
BatchNorm2d-21	[-1, 256, 8, 8]	512
Conv2d-22	[-1, 256, 8, 8]	589,824
BatchNorm2d-23	[-1, 256, 8, 8]	512
Conv2d-24	[-1, 256, 8, 8]	32,768
BatchNorm2d-25	[-1, 256, 8, 8]	512
BasicBlock-26	[-1, 256, 8, 8]	0
Conv2d-27	[-1, 512, 4, 4]	1,179,448
BatchNorm2d-28	[-1, 512, 4, 4]	1,024
Conv2d-29	[-1, 512, 4, 4]	2,359,296
BatchNorm2d-30	[-1, 512, 4, 4]	1,024
Conv2d-31	[-1, 512, 4, 4]	131,672
BatchNorm2d-32	[-1, 512, 4, 4]	1,024
BasicBlock-33	[-1, 512, 4, 4]	0
Linear-34	[-1, 10]	5,130
Total params: 4,977,226		
Trainable params: 4,977,226		
Non-trainable params: 0		
Input size (MB): 0.01		
Forward/backward pass size (MB): 9.06		
Params size (MB): 18.99		
Estimated Total Size (MB): 28.06		

Figure 5: Parameters

```

Epoch [42], train_loss: 0.0658, val_loss: 0.3301, val_acc: 0.9198
Epoch [43], train_loss: 0.0659, val_loss: 0.3447, val_acc: 0.9162
Epoch [44], train_loss: 0.0670, val_loss: 0.3456, val_acc: 0.9155
Epoch [45], train_loss: 0.0473, val_loss: 0.3411, val_acc: 0.9200
Epoch [46], train_loss: 0.0485, val_loss: 0.3276, val_acc: 0.9242
Epoch [47], train_loss: 0.0425, val_loss: 0.3478, val_acc: 0.9231
Epoch [48], train_loss: 0.0416, val_loss: 0.3399, val_acc: 0.9204
Epoch [49], train_loss: 0.0325, val_loss: 0.3344, val_acc: 0.9276
Epoch [50], train_loss: 0.0323, val_loss: 0.3415, val_acc: 0.9256
Epoch [51], train_loss: 0.0289, val_loss: 0.3422, val_acc: 0.9258
Epoch [52], train_loss: 0.0248, val_loss: 0.3405, val_acc: 0.9292
Epoch [53], train_loss: 0.0198, val_loss: 0.3615, val_acc: 0.9284
Epoch [54], train_loss: 0.0166, val_loss: 0.3590, val_acc: 0.9292
Epoch [55], train_loss: 0.0140, val_loss: 0.3625, val_acc: 0.9316
Epoch [56], train_loss: 0.0111, val_loss: 0.3652, val_acc: 0.9314
Epoch [57], train_loss: 0.0093, val_loss: 0.3672, val_acc: 0.9337
Epoch [58], train_loss: 0.0084, val_loss: 0.3591, val_acc: 0.9355
Epoch [59], train_loss: 0.0074, val_loss: 0.3599, val_acc: 0.9347
Epoch [60], train_loss: 0.0058, val_loss: 0.3446, val_acc: 0.9371
Epoch [61], train_loss: 0.0047, val_loss: 0.3543, val_acc: 0.9361
Epoch [62], train_loss: 0.0045, val_loss: 0.3544, val_acc: 0.9377
Epoch [63], train_loss: 0.0038, val_loss: 0.3455, val_acc: 0.9375
Epoch [64], train_loss: 0.0032, val_loss: 0.3466, val_acc: 0.9377
Epoch [65], train_loss: 0.0029, val_loss: 0.3511, val_acc: 0.9371
Epoch [66], train_loss: 0.0026, val_loss: 0.3462, val_acc: 0.9381
Epoch [67], train_loss: 0.0027, val_loss: 0.3475, val_acc: 0.9386
Epoch [68], train_loss: 0.0027, val_loss: 0.3521, val_acc: 0.9379
Epoch [69], train_loss: 0.0028, val_loss: 0.3492, val_acc: 0.9381
CPU times: user 35min 50s, sys: 6.56 s, total: 35min 56s
Wall time: 35min 45s

```

Figure 6: Parameters

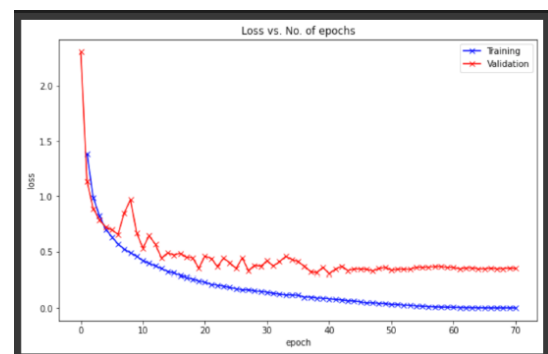


Figure 7: Loss vs. No. of epochs

that the learning rate starts at a low value, and gradually increases for 0.3 of the iterations to a maximum value of 0.01, and then gradually decreases to a very small value.

Findings and other testing results

Cutting the number of channel for each residual layers, the number of trainable parameters drops dramatically. Adding

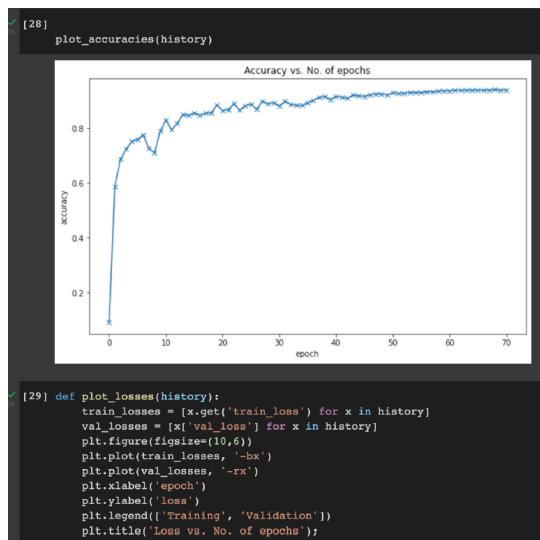


Figure 8: Accuracy vs No. of epochs

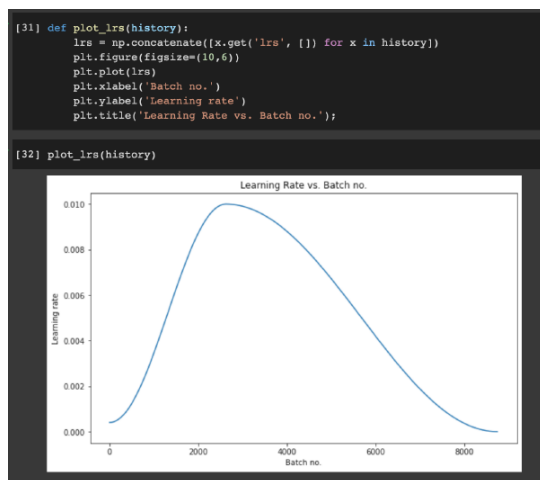


Figure 9: Learning Rates vs No. of epochs

more residual blocks in each layers the trainable parameters increases. For comparison, we keep two different model's trainable parameters around 4.9 million. The first model has original number of channel and one residual block in each layers. The second model cuts the number of channel in half for each layer and increases the residual block number to 6 for first layer, 5 for second layer, 4 for third layer and 3 for fourth layer. The best accuracy we get from first model is 0.9338 and the best accuracy we get from second model is 0.9387. Apparently the second model is slightly better than the first one. Therefore increasing the residual block number for each layer has better reward compared to increasing the channel numbers for each layer.

Our main findings are listed as follows:

- The default ResNet-18 model is excessive in CIFAR-10 classification. Simpler models are sufficient to satisfy this project goal. Our model could reach around an accuracy

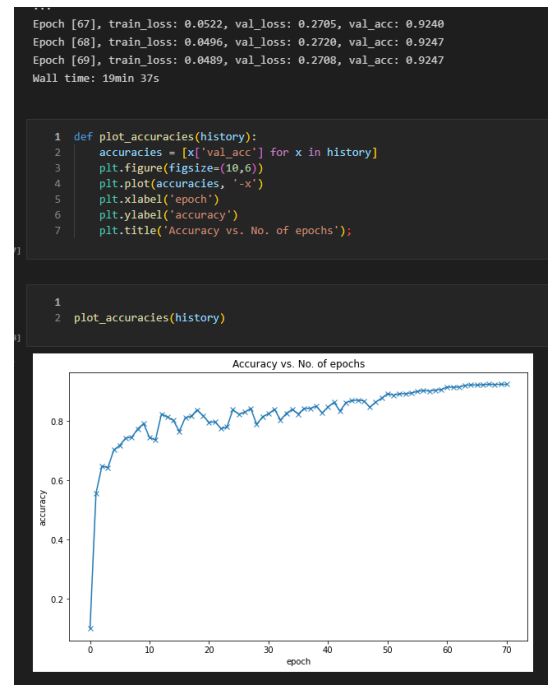


Figure 10: Accuracy vs No. of epochs

of almost accuracy 0.94 with far less than 5 million parameters.

- The width of ResNet design is more important than the depth. With holding the same parameter numbers, a deeper ResNet may work worse than a wider one. This will illustrate the idea of ResNet. As shown in the figure 10, a ResNet model with less than 1 million parameters could also reach more than 0.92 accuracy. It proved that it is more important to do modifications on the width of the model.
- In terms of data normalization and augmentation, it is useful. Normalization helps to get the data within a range (we specify) and which helps in making training a lot faster. Image augmentation is a super effective concept when we don't have enough data with us.
- In terms of hyper parameters, a too high learning rate leads to a suffering at first, large batch size makes little difference in our test.

References

He K, Zhang X, Ren S, et al. Deep residual learning for image recognition[C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2016: 770-778.