

# Quantization

## Quantize weight and input

To quantize the weight of the net, we can use `model.state_dict()` to get the weight and bias of the net, and modify the value of them. After the changes, we can use `model.load_state_dict(state)` to load the new weights and biases. But it is necessary to clone a new model, otherwise the training will be interrupted because of the loading data. The following is the related code:

```
1  model_q = copy.deepcopy(model)
2  state = model.state_dict()
3  state_dict_quant = OrderedDict()
4
5  for k, v in state.items():
6      v_quant = v
7
8      if 'features' in k:      # conv layer
9          v_quant = quant.linear_quantize(v, args.linear_bits)
10
11     if 'classifier' in k:    # linear layer
12         v_quant = quant.linear_quantize(v, args.conv_bits)
13
14     if 'running' in k:
15         state_dict_quant[k] = v
16         continue
17
18     if v.nelement() != 0:
19         v_quant = quant.linear_quantize(v, args.linear_bits)
20
21     state_dict_quant[k] = v_quant
22
23  model_q.load_state_dict(state_dict_quant)
```

But there is a problem stops me, the weight after training is mostly decimal. After convert it to `int8`, all of the weights turn to zero, and the inference can not work at all. To solve this problem, I firstly normalize the data to  $[0, 1]$ , then I expand the range to  $[-2^{bits-1}, 2^{bits-1} - 1]$ .

Quantization of a tensor is implemented in `linear_quantize`:

```
1  import torch
2  import math
3
4  def linear_quantize(input, bits):
5
6      range1 = input.max() - input.min()
7      temp = (input - input.min()) / range1
8
9      range2 = math.pow(2.0, bits) - 1
10     min_v = - math.pow(2.0, bits - 1)
11     ret = (temp.float() * range2) + min_v
12     ret = ret.float()
13
14     return ret
```

First, the input is scaled to  $[-2^{bits-1}, 2^{bits-1} - 1]$ . Then, use `int()` remove the decimal. After that, use `float()` change the to the data type which supported in GPU.

## Integer

In this part, I use the above `linear_quantize` function to handle all the weight and bias. As for input data, I process it using `input.int().float()`

- Run Docker

```
1 $ sudo nvidia-docker run -it --ipc=host --rm -v /home:/workspace
  nvcr.io/nvidia/pytorch:19.06-py3
2
3 $ cd leafz/PyTorch-Learning/quantization
```

- 16 bit

```
1 $ python -m torch.distributed.launch --nproc_per_node=4 resnet_q.py -a resnet50 --b 256 --
  epochs 50 --workers 4 --opt-level 02 --conv-bits 16 --linear-bits 16 ./
```

Final epoch: `Prec@1 63.850 Prec@5 93.800`

- 8 bit

```
1 $ python -m torch.distributed.launch --nproc_per_node=4 resnet_q.py -a resnet50 --b 256 --
  epochs 50 --workers 4 --opt-level 02 --conv-bits 8 --linear-bits 8 ./
```

Final epoch: `Prec@1 62.520 Prec@5 93.030`

- 4 bit

```
1 $ python -m torch.distributed.launch --nproc_per_node=4 resnet_q.py -a resnet50 --b 256 --
  epochs 50 --workers 4 --opt-level 02 --conv-bits 4 --linear-bits 4 ./
```

Final epoch: `Prec@1 64.170 Prec@5 93.820`

## Float

In this part, I remove the `int()` to get a result of control group in order to see how much influence decimal has on the above results.

- 16 bit

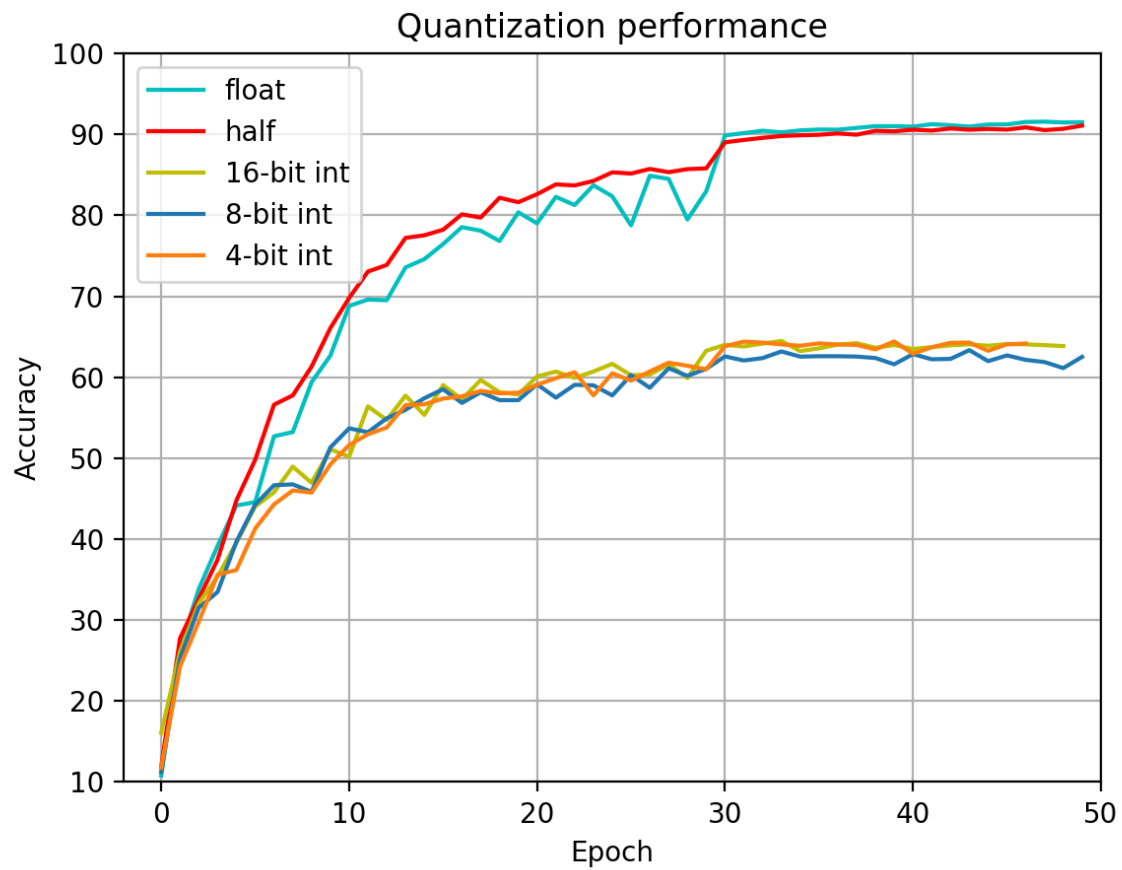
I change all the data type to the half in this part:

```
1 $ python -m torch.distributed.launch --nproc_per_node=4 resnet_q.py -a resnet50 --b 256 --  
epochs 50 --workers 4 --opt-level 02 --conv-bits 16 --linear-bits 16 ./
```

Final epoch: Prec@1 91.090 Prec@5 99.660

## Result

I plot the result of float, half, 4-bit int, 8-bit int and 16-bit int:



We can see that the accuracy of int are quite low. And the differences between several int are very small. I think the decimal should not influence so much.