

__vgg16__quant__ws__prune

December 3, 2024

```
[1]: import os
import time
import shutil

import torch
import torch.nn as nn

import torchvision
import torchvision.transforms as transforms

import torch.nn.utils.prune as prune

from models import *
from models.prune_util import *

import os
os.environ["CUDA_DEVICE_ORDER"]="PCI_BUS_ID"
os.environ["CUDA_VISIBLE_DEVICES"]="0"

import gc
gc.collect()
torch.cuda.empty_cache()

global best_prec

batch_size = 64
model_name = f"VGG16_ws_iter_prune_0.78_q"
fdir = 'result/' + model_name
model = VGG16()
ws_prune_vgg16(model, 0.78)
checkpoint = torch.load(f"result/VGG16_ws_iter_prune_0.78/model_best.pth.tar")
model.load_state_dict(checkpoint['state_dict'])
model.cuda()

device = torch.device("cuda")
```

```

normalize = transforms.Normalize(mean=[0.491, 0.482, 0.447], std=[0.247, 0.243,
↪0.262])

train_dataset = torchvision.datasets.CIFAR10(
    root='./data',
    train=True,
    download=True,
    transform=transforms.Compose([
        transforms.RandomCrop(32, padding=4),
        transforms.RandomHorizontalFlip(),
        transforms.ToTensor(),
        normalize,
    ]))
trainloader = torch.utils.data.DataLoader(train_dataset, batch_size=batch_size,
↪shuffle=True)

test_dataset = torchvision.datasets.CIFAR10(
    root='./data',
    train=False,
    download=True,
    transform=transforms.Compose([
        transforms.ToTensor(),
        normalize,
    ]))
testloader = torch.utils.data.DataLoader(test_dataset, batch_size=batch_size,
↪shuffle=False)

print_freq = 100 # every 100 batches, accuracy printed. Here, each batch
↪includes "batch_size" data points
# CIFAR10 has 50,000 training data, and 10,000 validation data.

def train(trainloader, model, criterion, optimizer, epoch):
    batch_time = AverageMeter()
    data_time = AverageMeter()
    losses = AverageMeter()
    top1 = AverageMeter()

    model.train()

    end = time.time()
    for i, (input, target) in enumerate(trainloader):
        # measure data loading time
        data_time.update(time.time() - end)

```

```

input, target = input.cuda(), target.cuda()

# compute output
output = model(input)
loss = criterion(output, target)

# measure accuracy and record loss
prec = accuracy(output, target)[0]
losses.update(loss.item(), input.size(0))
top1.update(prec.item(), input.size(0))

# compute gradient and do SGD step
optimizer.zero_grad()
loss.backward()
optimizer.step()

# measure elapsed time
batch_time.update(time.time() - end)
end = time.time()

if i % print_freq == 0:
    print('Epoch: [{0}] [{1}/{2}]\t'
          'Time {batch_time.val:.3f} ({batch_time.avg:.3f})\t'
          'Data {data_time.val:.3f} ({data_time.avg:.3f})\t'
          'Loss {loss.val:.4f} ({loss.avg:.4f})\t'
          'Prec {top1.val:.3f}% ({top1.avg:.3f}%)'.format(
            epoch, i, len(trainloader), batch_time=batch_time,
            data_time=data_time, loss=losses, top1=top1))

def validate(val_loader, model, criterion ):
    batch_time = AverageMeter()
    losses = AverageMeter()
    top1 = AverageMeter()

    # switch to evaluate mode
    model.eval()

    end = time.time()
    with torch.no_grad():
        for i, (input, target) in enumerate(val_loader):

            input, target = input.cuda(), target.cuda()

            # compute output
            output = model(input)
            loss = criterion(output, target)

```

```

        # measure accuracy and record loss
        prec = accuracy(output, target)[0]
        losses.update(loss.item(), input.size(0))
        top1.update(prec.item(), input.size(0))

        # measure elapsed time
        batch_time.update(time.time() - end)
        end = time.time()

        if i % print_freq == 0: # This line shows how frequently print out
            ↪ the status. e.g., i%5 => every 5 batch, prints out
            print('Test: [{0}/{1}]\t'
                  'Time {batch_time.val:.3f} ({batch_time.avg:.3f})\t'
                  'Loss {loss.val:.4f} ({loss.avg:.4f})\t'
                  'Prec {top1.val:.3f}% ({top1.avg:.3f}%)'.format(
                      i, len(val_loader), batch_time=batch_time, loss=losses,
                      top1=top1))

        print(' * Prec {top1.avg:.3f}% '.format(top1=top1))
        return top1.avg

def accuracy(output, target, topk=(1,)):
    """Computes the precision@k for the specified values of k"""
    maxk = max(topk)
    batch_size = target.size(0)

    _, pred = output.topk(maxk, 1, True, True)
    pred = pred.t()
    correct = pred.eq(target.view(1, -1).expand_as(pred))

    res = []
    for k in topk:
        correct_k = correct[:k].view(-1).float().sum(0)
        res.append(correct_k.mul_(100.0 / batch_size))
    return res

class AverageMeter(object):
    """Computes and stores the average and current value"""
    def __init__(self):
        self.reset()

    def reset(self):
        self.val = 0
        self.avg = 0

```

```

        self.sum = 0
        self.count = 0

    def update(self, val, n=1):
        self.val = val
        self.sum += val * n
        self.count += n
        self.avg = self.sum / self.count

def save_checkpoint(state, is_best, fdir):
    filepath = os.path.join(fdir, 'checkpoint.pth')
    torch.save(state, filepath)
    if is_best:
        shutil.copyfile(filepath, os.path.join(fdir, 'model_best.pth.tar'))

def adjust_learning_rate(optimizer, epoch):
    """For resnet, the lr starts from 0.1, and is divided by 10 at 80 and 120_
    epochs"""
    adjust_list = [150, 225]
    if epoch in adjust_list:
        for param_group in optimizer.param_groups:
            param_group['lr'] = param_group['lr'] * 0.1

def train_model(model, fdir, criterion, optimizer, epochs):
    os.makedirs(fdir, exist_ok=True)

    best_prec = 0

    #model = nn.DataParallel(model).cuda()
    model.cuda()
    criterion = criterion.cuda()
    #cudnn.benchmark = True

    for epoch in range(0, epochs):
        adjust_learning_rate(optimizer, epoch)

        train(trainloader, model, criterion, optimizer, epoch)

        # evaluate on test set
        print("Validation starts")
        prec = validate(testloader, model, criterion)

        # remember best precision and save checkpoint
        is_best = prec > best_prec

```

```

best_prec = max(prec, best_prec)
print('best acc: {:.1f}'.format(best_prec))
save_checkpoint({
    'epoch': epoch + 1,
    'state_dict': model.state_dict(),
    'best_prec': best_prec,
    'optimizer': optimizer.state_dict(),
}, is_best, fdir)

def val_model(model):

    model.cuda()
    model.eval()

    test_loss = 0
    correct = 0

    with torch.no_grad():
        for data, target in testloader:
            data, target = data.to(device), target.to(device) # loading to GPU
            output = model(data)
            pred = output.argmax(dim=1, keepdim=True)
            correct += pred.eq(target.view_as(pred)).sum().item()

    test_loss /= len(testloader.dataset)

    print('\nTest set: Accuracy: {}/{} ({:.0f}%) \n'.format(
        correct, len(testloader.dataset),
        100. * correct / len(testloader.dataset)))

```

```

Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)

```

/tmp/ipykernel_97379/3510108852.py:31: FutureWarning: You are using `torch.load` with `weights_only=False` (the current default value), which uses the default pickle module implicitly. It is possible to construct malicious pickle data which will execute arbitrary code during unpickling (See <https://github.com/pytorch/pytorch/blob/main/SECURITY.md#untrusted-models> for

more details). In a future release, the default value for `weights_only` will be flipped to `True`. This limits the functions that could be executed during unpickling. Arbitrary objects will no longer be allowed to be loaded via this mode unless they are explicitly allowlisted by the user via `torch.serialization.add_safe_globals`. We recommend you start setting `weights_only=True` for any use case where you don't have full control of the loaded file. Please open an issue on GitHub for any issues related to this experimental feature.

```
checkpoint = torch.load(f"result/VGG16_ws_iter_prune_0.78/model_best.pth.tar")
```

Files already downloaded and verified

Files already downloaded and verified

```
[2]: val_model(model)
```

Test set: Accuracy: 9077/10000 (91%)

```
[3]: quantize_pruned(model)
```

```
val_model(model)
```

Test set: Accuracy: 1322/10000 (13%)

```
[ ]: lr = 3e-3
```

```
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=lr)
train_model(model, fdir, criterion, optimizer, 100)
```

```
[4]: model = VGG16()
ws_prune_vgg16(model, 0.78)
quantize_pruned(model)

PATH = f"{fdir}/model_best.pth.tar"
checkpoint = torch.load(PATH)
model.load_state_dict(checkpoint['state_dict'])
model.cuda()
```

```
class SaveOutput:
    def __init__(self):
        self.outputs = []
    def __call__(self, module, module_in):
        self.outputs.append(module_in)
    def clear(self):
```

```

        self.outputs = []
save_output = SaveOutput()
model.features[40].register_forward_pre_hook(save_output)

val_model(model)

print_sparsity(model)

```

```

Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)
Pruning 7 kij-sticks out of 9 kij-sticks per output channel (77.8% pruned)

```

/tmp/ipykernel_97379/3671323971.py:6: FutureWarning: You are using `torch.load` with `weights_only=False` (the current default value), which uses the default pickle module implicitly. It is possible to construct malicious pickle data which will execute arbitrary code during unpickling (See <https://github.com/pytorch/pytorch/blob/main/SECURITY.md#untrusted-models> for more details). In a future release, the default value for `weights_only` will be flipped to `True`. This limits the functions that could be executed during unpickling. Arbitrary objects will no longer be allowed to be loaded via this mode unless they are explicitly allowlisted by the user via `torch.serialization.add_safe_globals`. We recommend you start setting `weights_only=True` for any use case where you don't have full control of the loaded file. Please open an issue on GitHub for any issues related to this experimental feature.

```
checkpoint = torch.load(PATH)
```

Test set: Accuracy: 9079/10000 (91%)

```

layer 3 sparsity: 0.778
layer 7 sparsity: 0.778
layer 10 sparsity: 0.778
layer 14 sparsity: 0.778
layer 17 sparsity: 0.778
layer 20 sparsity: 0.778
layer 24 sparsity: 0.778
layer 27 sparsity: 0.778
layer 30 sparsity: 0.778

```



```
layer 34 sparsity: 0.778
layer 37 sparsity: 0.778
layer 40 sparsity: 0.778
```

```
[16]: print(model.features[40].weight_q.data[:2,:2])
```

```
tensor([[[[ 0.0000,  5.1863,  0.0000],
           [ 0.0000,  1.2966,  0.0000],
           [ 0.0000,  0.0000,  0.0000]],

         [[ 0.0000,  0.0000,  0.0000],
           [ 0.0000, -1.2966,  0.0000],
           [ 0.0000,  0.0000,  0.0000]]],

        [[[ 0.0000,  0.0000,  0.0000],
           [ 0.0000,  0.0000,  0.0000],
           [ 0.0000, -0.0000, -2.5932]],

         [[ 0.0000,  0.0000,  0.0000],
           [ 0.0000,  0.0000,  0.0000],
           [ 0.0000, -0.0000, -1.2966]]]], device='cuda:0')
```

```
[21]: print(save_output.outputs[0][0][:2,:2])
```

```
tensor([[[[0.0000, 0.0000],
           [0.0062, 0.0000]],

         [[0.0000, 0.0000],
           [0.0000, 0.0000]]],

        [[[0.6084, 0.0000],
           [0.0062, 1.4861]],

         [[0.0000, 0.0000],
           [0.0000, 0.0000]]]], device='cuda:0')
```

```
[14]: print(f'Weight int: \n{(model.features[40].weight_q.data / (model.features[40].
    ↪weight_quant.wgt_alpha.data.item()/(2**(4-1)-1))) [0,0]}')
x = save_output.outputs[0][0]
x_alpha = model.features[40].act_alpha.data.item()
x_delta = x_alpha / (2**(4)-1)
act_q = act_quantization(4)
x_q = act_q(x, x_alpha)
print(f'Act int: \n{(x_q/x_delta)[:2,:2]}')
```

Weight int:

```

tensor([[0., 4., 0.],
        [0., 1., 0.],
        [0., 0., 0.]], device='cuda:0')
Act int:
tensor([[[[0.0000, 0.0000],
          [0.0000, 0.0000]],

         [[0.0000, 0.0000],
          [0.0000, 0.0000]]],

        [[[1.0000, 0.0000],
          [0.0000, 2.0000]],

         [[0.0000, 0.0000],
          [0.0000, 0.0000]]]], device='cuda:0')

```