

CS 182 Homework 10: Pruning

This notebook has been adapted with permission from MIT 6.S965 Fall 2022. Original authors: Yujun Lin, Ji Lin, Zhijian Liu and Song Han

Goals

In this assignment, you will practice pruning a classical neural network model to reduce both model size and latency. The goals of this assignment are as follows:

- Understand the basic concept of **pruning**
- Implement and apply **fine-grained pruning**
- Implement and apply **channel pruning**
- Get a basic understanding of performance improvement (such as speedup) from pruning
- Understand the differences and tradeoffs between these pruning approaches

Contents

There are two main sections in this lab: ***Fine-grained Pruning*** and ***Channel Pruning***. Questions 6-9 are OPTIONAL.

There are **9** questions in total:

- For *Fine-grained Pruning*, there are **5** questions (Question 1-5).
- For *Channel Pruning*, there are **3** questions (Question 6-8).
- Question 9 compares fine-grained pruning and channel pruning.

Setup

First, install the required packages and download the datasets and pretrained model. Here we use CIFAR10 dataset and VGG network which is the same as what we used in the Lab 0 tutorial.

```
In [1]: print('Installing torchprofile...')
#!pip install torchprofile 1>/dev/null
print('All required packages have been successfully installed!')
```

```
Installing torchprofile...
All required packages have been successfully installed!
```

```
In [2]: import os
os.environ["KMP_DUPLICATE_LIB_OK"]="TRUE"
```

```
In [3]: import copy
import math
import random
import time
from collections import OrderedDict, defaultdict
from typing import Union, List

import numpy as np
import torch
from matplotlib import pyplot as plt
from torch import nn
from torch.optim import *
from torch.optim.lr_scheduler import *
from torch.utils.data import DataLoader
from torchprofile import profile_macs
from torchvision.datasets import *
from torchvision.transforms import *
from tqdm.auto import tqdm

from torchprofile import profile_macs

assert torch.cuda.is_available(), \
"The current runtime does not have CUDA support." \
"Please go to menu bar (Runtime - Change runtime type) and select GPU"
```

```
In [4]: random.seed(0)
np.random.seed(0)
torch.manual_seed(0)
```

Out[4]: <torch._C.Generator at 0x1943dd5d690>

```
In [5]: def download_url(url, model_dir='.', overwrite=False):
import os, sys
from urllib.request import urlretrieve
target_dir = url.split('/')[-1]
model_dir = os.path.expanduser(model_dir)
try:
    if not os.path.exists(model_dir):
        os.makedirs(model_dir)
    model_dir = os.path.join(model_dir, target_dir)
    cached_file = model_dir
    if not os.path.exists(cached_file) or overwrite:
        sys.stderr.write('Downloading: "{}" to {}\n'.format(url, cached_file))
        urlretrieve(url, cached_file)
    return cached_file
except Exception as e:
    # remove lock file so download can be executed next time.
    os.remove(os.path.join(model_dir, 'download.lock'))
    sys.stderr.write('Failed to download from url %s' % url + '\n' + str(e) + '\n')
    return None
```

```

In [6]: class VGG(nn.Module):
    ARCH = [64, 128, 'M', 256, 256, 'M', 512, 512, 'M', 512, 512, 'M']

    def __init__(self) -> None:
        super().__init__()

        layers = []
        counts = defaultdict(int)

        def add(name: str, layer: nn.Module) -> None:
            layers.append((f"{name} {counts[name]}", layer))
            counts[name] += 1

        in_channels = 3
        for x in self.ARCH:
            if x != 'M':
                # conv-bn-relu
                add("conv", nn.Conv2d(in_channels, x, 3, padding=1, bias=False))
                add("bn", nn.BatchNorm2d(x))
                add("relu", nn.ReLU(True))
                in_channels = x
            else:
                # maxpool
                add("pool", nn.MaxPool2d(2))

        self.backbone = nn.Sequential(OrderedDict(layers))
        self.classifier = nn.Linear(512, 10)

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        # backbone: [N, 3, 32, 32] => [N, 512, 2, 2]
        x = self.backbone(x)

        # avgpool: [N, 512, 2, 2] => [N, 512]
        x = x.mean([2, 3])

        # classifier: [N, 512] => [N, 10]
        x = self.classifier(x)
        return x

```

```
In [7]: def train(
        model: nn.Module,
        dataloader: DataLoader,
        criterion: nn.Module,
        optimizer: Optimizer,
        scheduler: LambdaLR,
        callbacks = None
    ) -> None:
    model.train()

    for inputs, targets in tqdm(dataloader, desc='train', leave=False):
        # Move the data from CPU to GPU
        inputs = inputs.cuda()
        targets = targets.cuda()

        # Reset the gradients (from the last iteration)
        optimizer.zero_grad()

        # Forward inference
        outputs = model(inputs)
        loss = criterion(outputs, targets)

        # Backward propagation
        loss.backward()

        # Update optimizer and LR scheduler
        optimizer.step()
        scheduler.step()

    if callbacks is not None:
        for callback in callbacks:
            callback()
```

```
In [8]: @torch.inference_mode()
def evaluate(
    model: nn.Module,
    dataloader: DataLoader,
    verbose=True,
) -> float:
    model.eval()

    num_samples = 0
    num_correct = 0

    for inputs, targets in tqdm(dataloader, desc="eval", leave=False,
                                disable=not verbose):
        # Move the data from CPU to GPU
        inputs = inputs.cuda()
        targets = targets.cuda()

        # Inference
        outputs = model(inputs)

        # Convert logits to class indices
        outputs = outputs.argmax(dim=1)

        # Update metrics
        num_samples += targets.size(0)
        num_correct += (outputs == targets).sum()

    return (num_correct / num_samples * 100).item()
```

Helper Functions (Flops, Model Size calculation, etc.)

```

In [9]: def get_model_macs(model, inputs) -> int:
        return profile_macs(model, inputs)

def get_sparsity(tensor: torch.Tensor) -> float:
    """
    calculate the sparsity of the given tensor
    sparsity = #zeros / #elements = 1 - #nonzeros / #elements
    """
    return 1 - float(tensor.count_nonzero()) / tensor.numel()

def get_model_sparsity(model: nn.Module) -> float:
    """
    calculate the sparsity of the given model
    sparsity = #zeros / #elements = 1 - #nonzeros / #elements
    """
    num_nonzeros, num_elements = 0, 0
    for param in model.parameters():
        num_nonzeros += param.count_nonzero()
        num_elements += param.numel()
    return 1 - float(num_nonzeros) / num_elements

def get_num_parameters(model: nn.Module, count_nonzero_only=False) -> int:
    """
    calculate the total number of parameters of model
    :param count_nonzero_only: only count nonzero weights
    """
    num_counted_elements = 0
    for param in model.parameters():
        if count_nonzero_only:
            num_counted_elements += param.count_nonzero()
        else:
            num_counted_elements += param.numel()
    return num_counted_elements

def get_model_size(model: nn.Module, data_width=32, count_nonzero_only=False) -> int:
    """
    calculate the model size in bits
    :param data_width: #bits per element
    :param count_nonzero_only: only count nonzero weights
    """
    return get_num_parameters(model, count_nonzero_only) * data_width

Byte = 8
KiB = 1024 * Byte
MiB = 1024 * KiB
GiB = 1024 * MiB

```

Define misc functions for verification.

```

In [10]: def test_fine_grained_prune(
    test_tensor=torch.tensor([[-0.46, -0.40, 0.39, 0.19, 0.37],
                                [0.00, 0.40, 0.17, -0.15, 0.16],
                                [-0.20, -0.23, 0.36, 0.25, 0.03],
                                [0.24, 0.41, 0.07, 0.13, -0.15],
                                [0.48, -0.09, -0.36, 0.12, 0.45]]),
    test_mask=torch.tensor([[True, True, False, False, False],
                             [False, True, False, False, False],
                             [False, False, False, False, False],
                             [False, True, False, False, False],
                             [True, False, False, False, True]]),
    target_sparsity=0.75, target_nonzeros=None):
    def plot_matrix(tensor, ax, title):
        ax.imshow(tensor.cpu().numpy() == 0, vmin=0, vmax=1, cmap='tab20c')
        ax.set_title(title)
        ax.set_yticklabels([])
        ax.set_xticklabels([])
        for i in range(tensor.shape[1]):
            for j in range(tensor.shape[0]):
                text = ax.text(j, i, f'{tensor[i, j].item():.2f}',
                               ha="center", va="center", color="k")

    test_tensor = test_tensor.clone()
    fig, axes = plt.subplots(1, 2, figsize=(6, 10))
    ax_left, ax_right = axes.ravel()
    plot_matrix(test_tensor, ax_left, 'dense tensor')

    sparsity_before_pruning = get_sparsity(test_tensor)
    mask = fine_grained_prune(test_tensor, target_sparsity)
    sparsity_after_pruning = get_sparsity(test_tensor)
    sparsity_of_mask = get_sparsity(mask)

    plot_matrix(test_tensor, ax_right, 'sparse tensor')
    fig.tight_layout()
    plt.show()

    print('* Test fine_grained_prune()')
    print(f'    target sparsity: {target_sparsity:.2f}')
    print(f'    sparsity before pruning: {sparsity_before_pruning:.2f}')
    print(f'    sparsity after pruning: {sparsity_after_pruning:.2f}')
    print(f'    sparsity of pruning mask: {sparsity_of_mask:.2f}')

    if target_nonzeros is None:
        if test_mask.equal(mask):
            print('* Test passed.')
        else:
            print('* Test failed.')
    else:
        if mask.count_nonzero() == target_nonzeros:
            print('* Test passed.')
        else:
            print('* Test failed.')

```

Load the pretrained model and the CIFAR-10 dataset.

```
In [11]: #checkpoint_url = "https://hanlab.mit.edu/files/course/labs/vgg.cifar.pretrained.pth
checkpoint_url = "https://hanlab18.mit.edu/files/course/labs/vgg.cifar.pretrained.pt
checkpoint = torch.load(download_url(checkpoint_url), map_location="cpu")
print("download successfully")
model = VGG().cuda()
print(f"=> loading checkpoint '{checkpoint_url}'")
model.load_state_dict(checkpoint['state_dict'])
recover_model = lambda: model.load_state_dict(checkpoint['state_dict'])
```

download successfully
=> loading checkpoint 'https://hanlab18.mit.edu/files/course/labs/vgg.cifar.pretrained.pth'

```
In [12]: image_size = 32
transforms = {
    "train": Compose([
        RandomCrop(image_size, padding=4),
        RandomHorizontalFlip(),
        ToTensor(),
    ]),
    "test": ToTensor(),
}
dataset = {}
for split in ["train", "test"]:
    dataset[split] = CIFAR10(
        #root="data/cifar10",
        root="../cifar-10/cifar-10-batches-py/",
        train=(split == "train"),
        download=True,
        transform=transforms[split],
    )
dataloader = {}
for split in ['train', 'test']:
    dataloader[split] = DataLoader(
        dataset[split],
        batch_size=512,
        shuffle=(split == 'train'),
        num_workers=0,
        pin_memory=True,
    )
```

Files already downloaded and verified
Files already downloaded and verified

Let's First Evaluate the Accuracy and Model Size of Dense Model

Neural networks have become ubiquitous in many applications. Here we have loaded a pretrained VGG model for classifying images in CIFAR10 dataset.

Let's first evaluate the accuracy and model size of this model.


```
In [13]: dense_model_accuracy = evaluate(model, dataloader['test'])
dense_model_size = get_model_size(model)
print(f"dense model has accuracy={dense_model_accuracy:.2f}%")
print(f"dense model has size={dense_model_size/MiB:.2f} MiB")
```

```
eval: 0% | 0/20 [00:00<?, ?it/s]
```

```
dense model has accuracy=92.95%
```

```
dense model has size=35.20 MiB
```

While large neural networks are very powerful, their size consumes considerable storage, memory bandwidth, and computational resources. As we can see from the results above, a model for the task as simple as classifying 32×32 images into 10 classes can be as large as 35 MiB. For embedded mobile applications, these resource demands become prohibitive.

Therefore, neural network pruning is exploited to facilitates storage and transmission of mobile applications incorporating DNNs.

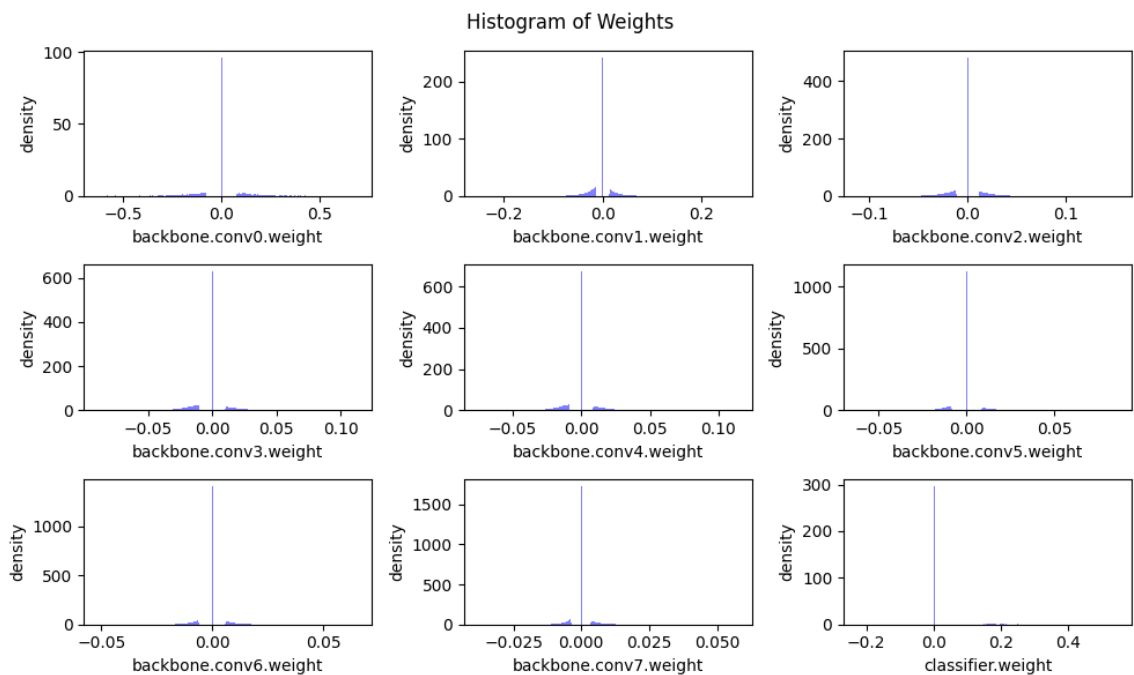
The goal of pruning is to reduce the model size while maintaining the accuracy.

Let's see the distribution of weight values

Before we jump into pruning, let's see the distribution of weight values in the dense model.

```
In [27]: def plot_weight_distribution(model, bins=256, count_nonzero_only=False):
fig, axes = plt.subplots(3,3, figsize=(10, 6))
axes = axes.ravel()
plot_index = 0
for name, param in model.named_parameters():
    if param.dim() > 1:
        ax = axes[plot_index]
        if count_nonzero_only:
            param_cpu = param.detach().view(-1).cpu()
            param_cpu = param_cpu[param_cpu != 0].view(-1)
            ax.hist(param_cpu, bins=bins, density=True,
                    color = 'blue', alpha = 0.5)
        else:
            ax.hist(param.detach().view(-1).cpu(), bins=bins, density=True,
                    color = 'blue', alpha = 0.5)
        ax.set_xlabel(name)
        ax.set_ylabel('density')
        plot_index += 1
fig.suptitle('Histogram of Weights')
fig.tight_layout()
fig.subplots_adjust(top=0.925)
plt.show()

plot_weight_distribution(model)
```



Question 1

Please answer the following questions using the information in the above histograms of weights.

Question 1.1

What are the common characteristics of the weight distribution in the different layers?

Your Answer:

Question 1.2

How do these characteristics help pruning?

Your Answer:

Fine-grained Pruning

In this section, we will implement and perform fine-grained pruning.

Fine-grained pruning removes the synapses with lowest importance. The weight tensor W will become sparse after fine-grained pruning, which can be described with **sparsity**:

$$\text{sparsity} := \#Zeros/\#W = 1 - \#Nonzeros/\#W$$

where $\#W$ is the number of elements in W .

In practice, given the target sparsity s , the weight tensor W is multiplied with a binary mask M to disregard removed weight:

$$\begin{aligned} v_{\text{thr}} &= \text{kthvalue}(\text{Importance}, \#W \cdot s) \\ M &= \text{Importance} > v_{\text{thr}} \\ W &= W \cdot M \end{aligned}$$

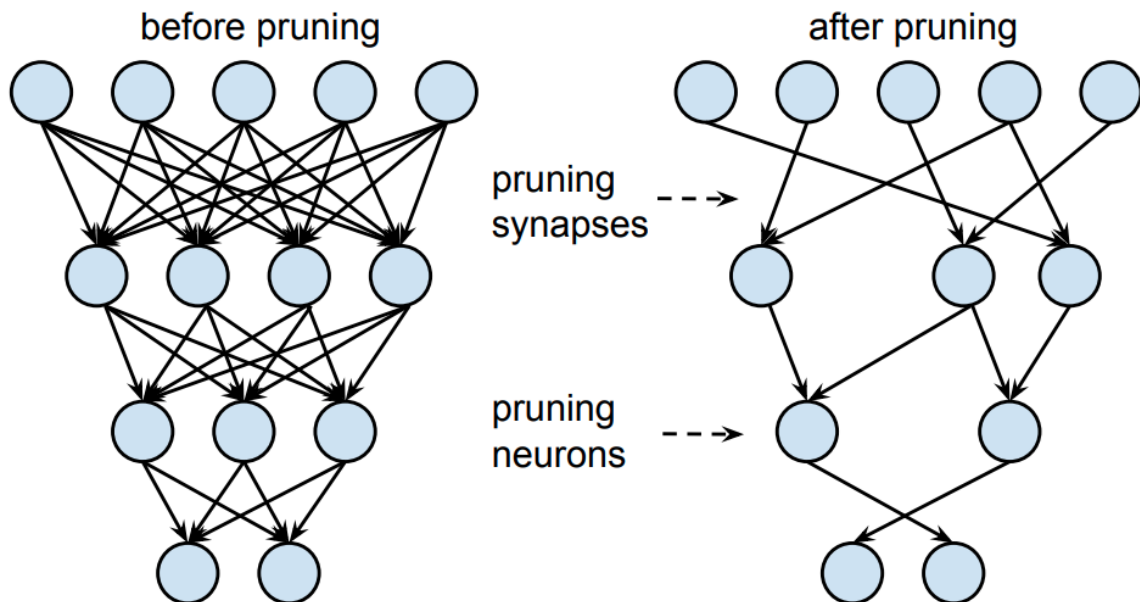
where Importance is importance tensor with the same shape of W , $\text{kthvalue}(X, k)$ finds the k -th smallest value of tensor X , v_{thr} is the threshold value.

Magnitude-based Pruning

For fine-grained pruning, a widely-used importance is the magnitude of weight value, *i.e.*,

$$\text{Importance} = |W|$$

This is known as **Magnitude-based Pruning** (see [Learning both Weights and Connections for Efficient Neural Networks \(https://arxiv.org/pdf/1506.02626.pdf\)](https://arxiv.org/pdf/1506.02626.pdf)).



Question 2

Please complete the following magnitude-based fine-grained pruning function.

Hint:

- In step 1, we calculate the number of zeros (`num_zeros`) after pruning. Note that `num_zeros` should be an integer. You could use either `round()` or `int()` to convert a floating number into an integer. Here we use `round()`.
- In step 2, we calculate the `importance` of weight tensor. Pytorch provides `torch.abs()` (<https://pytorch.org/docs/stable/generated/torch.abs.html#torch.abs>), `torch.Tensor.abs()` (<https://pytorch.org/docs/stable/generated/torch.Tensor.abs.html#torch.Tensor.abs>), `torch.Tensor.abs_()` (https://pytorch.org/docs/stable/generated/torch.Tensor.abs_.html) APIs.
- In step 3, we calculate the pruning `threshold` so that all synapses with importance smaller than `threshold` will be removed. Pytorch provides `torch.kthvalue()` (<https://pytorch.org/docs/stable/generated/torch.kthvalue.html>), `torch.Tensor.kthvalue()` (<https://pytorch.org/docs/stable/generated/torch.Tensor.kthvalue.html>), `torch.topk()` (<https://pytorch.org/docs/stable/generated/torch.topk.html>) APIs.
- In step 4, we calculate the pruning `mask` based on the `threshold`. `1` in the `mask` indicates the synapse will be kept, and `0` in the `mask` indicates the synapse will be removed. `mask = importance > threshold`. Pytorch provides `torch.gt()` (<https://pytorch.org/docs/stable/generated/torch.gt.html?highlight=torch%20gt#torch.gt>) API.

```

In [16]: def fine_grained_prune(tensor: torch.Tensor, sparsity : float) -> torch.Tensor:
        """
        magnitude-based pruning for single tensor
        :param tensor: torch.(cuda.)Tensor, weight of conv/fc layer
        :param sparsity: float, pruning sparsity
            sparsity = #zeros / #elements = 1 - #nonzeros / #elements
        :return:
            torch.(cuda.)Tensor, mask for zeros
        """

        sparsity = min(max(0.0, sparsity), 1.0)
        if sparsity == 1.0:
            tensor.zero_()
            return torch.zeros_like(tensor)
        elif sparsity == 0.0:
            return torch.ones_like(tensor)

        num_elements = tensor.numel()

        ##### YOUR CODE STARTS HERE #####
        # Step 1: calculate the #zeros (please use round())
        num_elements = tensor.numel()
        num_zeros = round(num_elements * sparsity)
        # Step 2: calculate the importance of weight
        importance = torch.abs(tensor)
        # Step 3: calculate the pruning threshold
        threshold = torch.kthvalue(importance.view(-1), k=num_zeros).values
        #print(f"threshold is {threshold}")
        # Step 4: get binary mask (1 for nonzeros, 0 for zeros)
        mask = torch.gt(importance, threshold)
        ##### YOUR CODE ENDS HERE #####

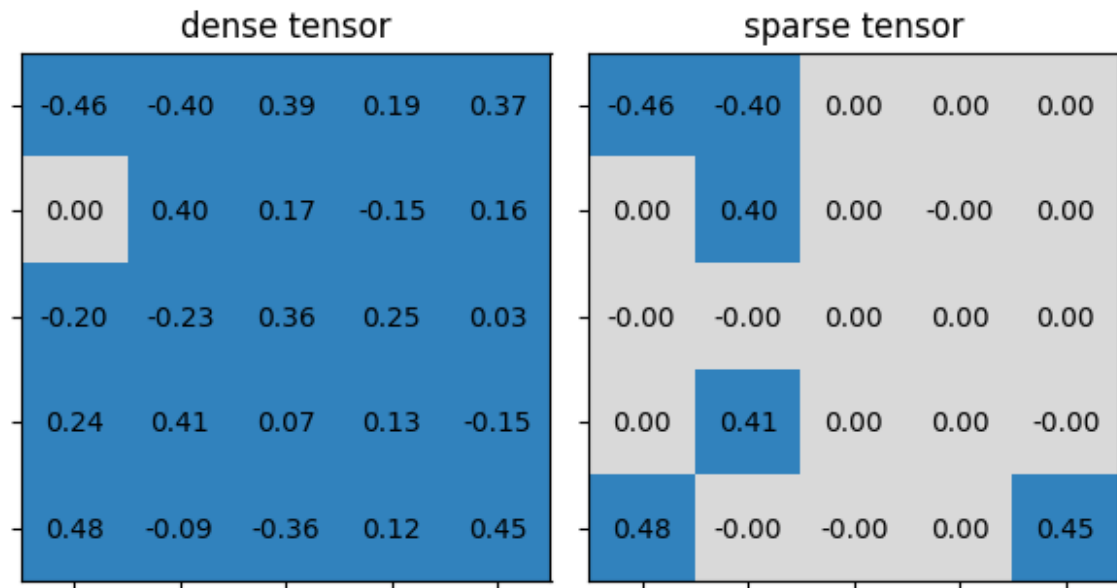
        # Step 5: apply mask to prune the tensor
        tensor.mul_(mask)

        return mask

```

Let's verify the functionality of defined fine-grained pruning by applying the function above on a dummy tensor.

```
In [17]: test_fine_grained_prune()
```



```
* Test fine_grained_prune()
  target sparsity: 0.75
    sparsity before pruning: 0.04
    sparsity after pruning: 0.76
    sparsity of pruning mask: 0.76
* Test passed.
```

Question 3

The last cell plots the tensor before and after pruning. Nonzeros are rendered in blue while zeros are rendered in gray. Please modify the value of `target_sparsity` in the following code cell so that there are only 10 nonzeros in the sparse tensor after pruning.

```
In [19]: ##### YOUR CODE STARTS HERE #####
# 10 / 25 = 0.4, so target_sparsity should be 1 - 0.4 = 0.6
target_sparsity = 0.6 # please modify the value of target_sparsity
##### YOUR CODE ENDS HERE #####
test_fine_grained_prune(target_sparsity=target_sparsity, target_nonzeros=10)
```



```
* Test fine_grained_prune()
  target sparsity: 0.60
    sparsity before pruning: 0.04
    sparsity after pruning: 0.60
    sparsity of pruning mask: 0.60
* Test passed.
```

We now wrap the fine-grained pruning function into a class for pruning the whole model. In class `FineGrainedPruner`, we have to keep a record of the pruning masks so that we could apply the masks whenever the model weights change to make sure the model keep sparse all the time.

```
In [20]: class FineGrainedPruner:
  def __init__(self, model, sparsity_dict):
    self.masks = FineGrainedPruner.prune(model, sparsity_dict)

  @torch.no_grad()
  def apply(self, model):
    for name, param in model.named_parameters():
      if name in self.masks:
        param *= self.masks[name]

  @staticmethod
  @torch.no_grad()
  def prune(model, sparsity_dict):
    masks = dict()
    for name, param in model.named_parameters():
      if param.dim() > 1: # we only prune conv and fc weights
        masks[name] = fine_grained_prune(param, sparsity_dict[name])
    return masks
```

Sensitivity Scan

Different layers contribute differently to the model performance. It is challenging to decide the proper sparsity for each layer. A widely used approach is sensitivity scan.

During the sensitivity scan, at each time, we will only prune one layer to see the accuracy degradation. By scanning different sparsities, we could draw the sensitivity curve (i.e., accuracy vs. sparsity) of the corresponding layer.

Here is an example figure for sensitivity curves. The x-axis is the sparsity or the percentage of #parameters dropped (i.e., sparsity). The y-axis is the validation accuracy. (This is Figure 6 in [Learning both Weights and Connections for Efficient Neural Networks](https://arxiv.org/pdf/1506.02626.pdf) (<https://arxiv.org/pdf/1506.02626.pdf>))

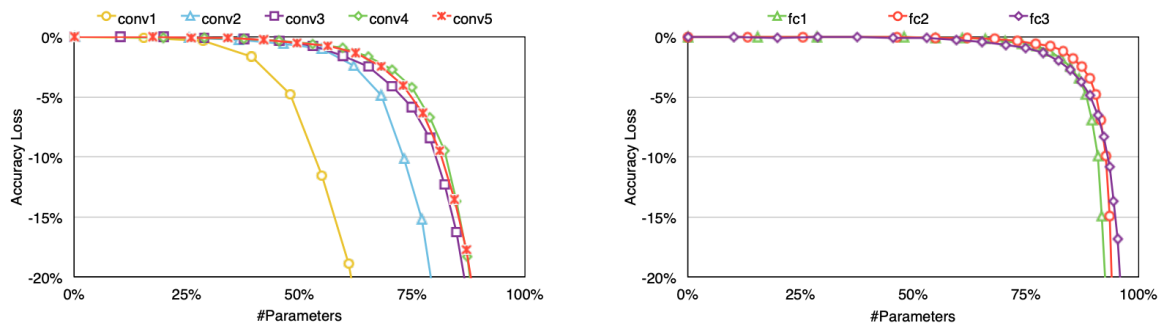


Figure 6: Pruning sensitivity for CONV layer (left) and FC layer (right) of AlexNet.

The following code cell defines the sensitivity scan function that returns the sparsities scanned, and a list of accuracies for each weight to be pruned.

```
In [21]: @torch.no_grad()
def sensitivity_scan(model, dataloader, scan_step=0.1, scan_start=0.4, scan_end=1.0):
    sparsities = np.arange(start=scan_start, stop=scan_end, step=scan_step)
    accuracies = []
    named_conv_weights = [(name, param) for (name, param) \
                           in model.named_parameters() if param.dim() > 1]
    for i_layer, (name, param) in enumerate(named_conv_weights):
        param_clone = param.detach().clone()
        accuracy = []
        for sparsity in tqdm(sparsities, desc=f'scanning {i_layer}/{len(named_conv_weights)}'):
            fine_grained_prune(param.detach(), sparsity=sparsity)
            acc = evaluate(model, dataloader, verbose=False)
            if verbose:
                print(f'\r    sparsity={sparsity:.2f}: accuracy={acc:.2f}%', end='')
            # restore
            param.copy_(param_clone)
            accuracy.append(acc)
        if verbose:
            print(f'\r    sparsity=[{",".join(["{: .2f} ".format(x) for x in sparsities]}]')
            accuracies.append(accuracy)
    return sparsities, accuracies
```

Please run the following cells to plot the sensitivity curves. It should take around 2 minutes to


```
In [22]: sparsities, accuracies = sensitivity_scan(
        model, dataloader['test'], scan_step=0.1, scan_start=0.4, scan_end=1.0)

scanning 0/9 weight - backbone.conv0.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.42%, 91.19%, 87.55%, 83.39%, 69.43%, 31.82%]

scanning 1/9 weight - backbone.conv1.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.93%, 92.88%, 92.71%, 92.40%, 91.32%, 84.78%]

scanning 2/9 weight - backbone.conv2.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.94%, 92.64%, 92.46%, 91.77%, 89.85%, 78.56%]

scanning 3/9 weight - backbone.conv3.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.86%, 92.72%, 92.23%, 91.09%, 85.35%, 51.29%]

scanning 4/9 weight - backbone.conv4.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.88%, 92.68%, 92.22%, 91.47%, 85.35%, 51.29%]

scanning 5/9 weight - backbone.conv5.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.92%, 92.71%, 92.64%, 91.88%, 89.90%, 82.21%]

scanning 6/9 weight - backbone.conv6.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.95%, 92.86%, 92.65%, 92.10%, 90.58%, 83.64%]

scanning 7/9 weight - backbone.conv7.weight: 0%|          | 0/6 [00:00<?, ?it/s]

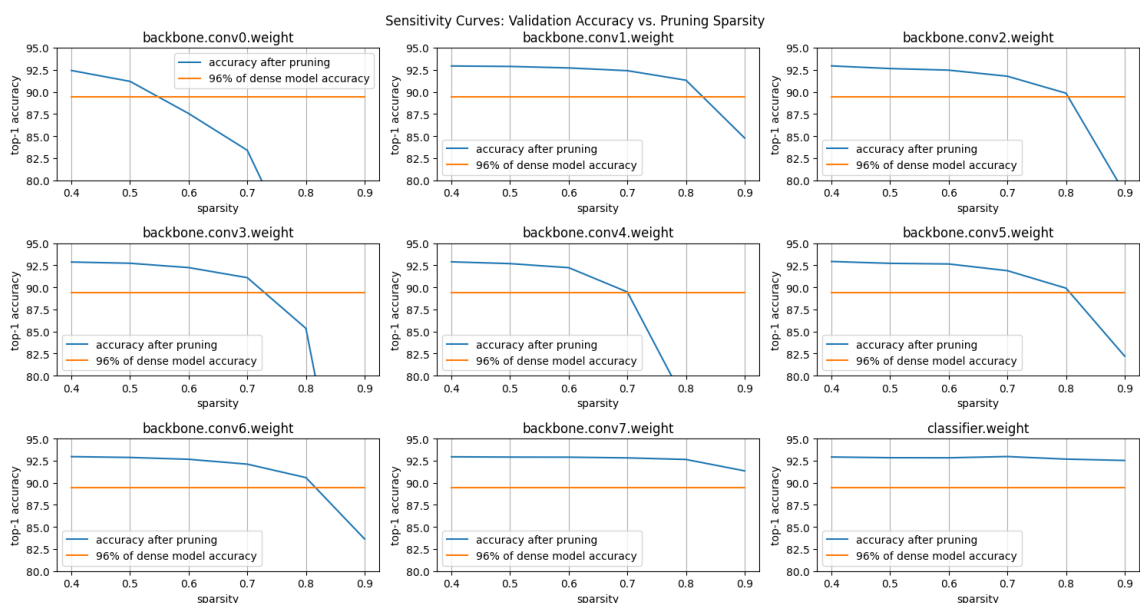
    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.93%, 92.90%, 92.89%, 92.81%, 92.63%, 91.34%]

scanning 8/9 weight - classifier.weight: 0%|          | 0/6 [00:00<?, ?it/s]

    sparsity=[0.40, 0.50, 0.60, 0.70, 0.80, 0.90]: accuracy=[92.91%, 92.83%, 92.82%, 92.96%, 92.67%, 92.52%]
```

```
In [23]: def plot_sensitivity_scan(sparsities, accuracies, dense_model_accuracy):
    lower_bound_accuracy = 100 - (100 - dense_model_accuracy) * 1.5
    fig, axes = plt.subplots(3, int(math.ceil(len(accuracies) / 3)), figsize=(15,8))
    axes = axes.ravel()
    plot_index = 0
    for name, param in model.named_parameters():
        if param.dim() > 1:
            ax = axes[plot_index]
            curve = ax.plot(sparsities, accuracies[plot_index])
            line = ax.plot(sparsities, [lower_bound_accuracy] * len(sparsities))
            ax.set_xticks(np.arange(start=0.4, stop=1.0, step=0.1))
            ax.set_ylim(80, 95)
            ax.set_title(name)
            ax.set_xlabel('sparsity')
            ax.set_ylabel('top-1 accuracy')
            ax.legend([
                'accuracy after pruning',
                f'{lower_bound_accuracy / dense_model_accuracy * 100:.0f}% of dense'
            ])
            ax.grid(axis='x')
            plot_index += 1
    fig.suptitle('Sensitivity Curves: Validation Accuracy vs. Pruning Sparsity')
    fig.tight_layout()
    fig.subplots_adjust(top=0.925)
    plt.show()

plot_sensitivity_scan(sparsities, accuracies, dense_model_accuracy)
```



Question 4

Please answer the following questions using the information in the above sensitivity curves.

Question 4.1

What's the relationship between pruning sparsity and model accuracy? (*i.e.*, does accuracy increase or decrease when sparsity becomes higher?)

Your Answer:

Question 4.2

Do all the layers have the same sensitivity?

Your Answer:

Question 4.3

Which layer is the most sensitive to the pruning sparsity?

Your Answer:

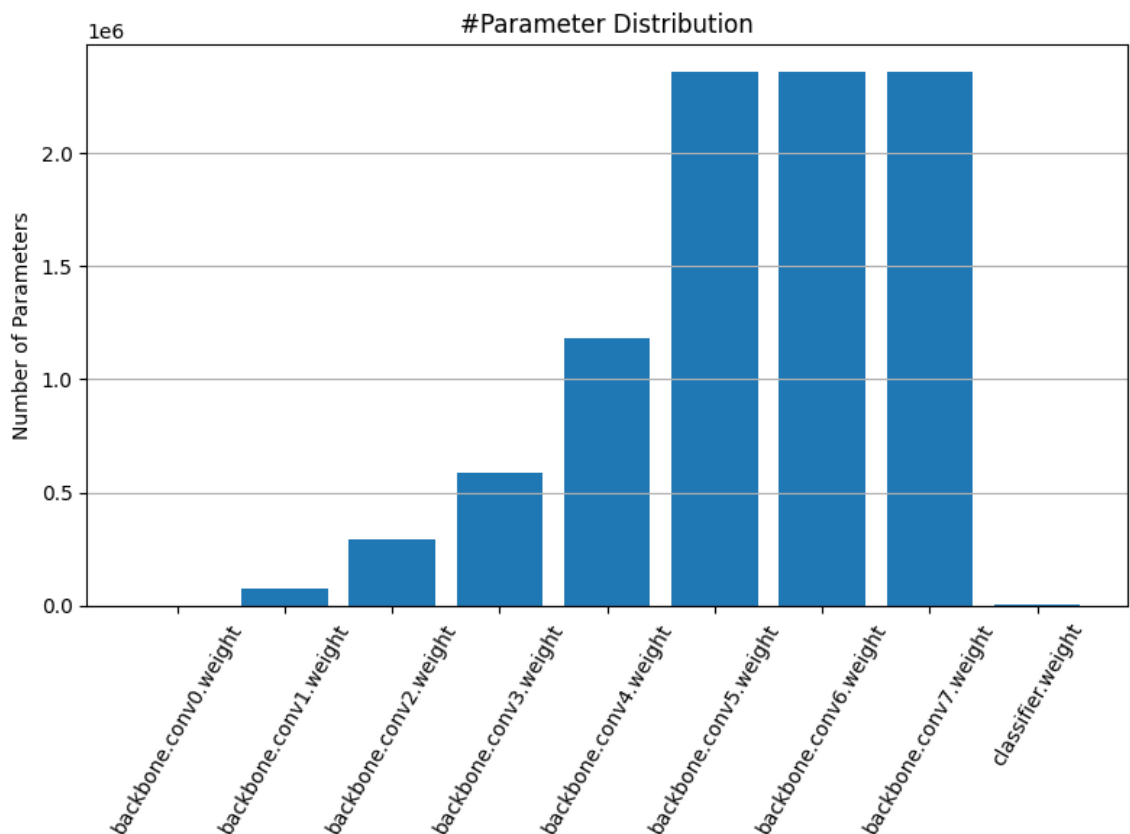
#Parameters of each layer

In addition to accuracy, the number of each layer's parameters also affects the decision on sparsity selection. Layers with more #parameters require larger sparsities.

Please run the following code cell to plot the distribution of #parameters in the whole model.

```
In [24]: def plot_num_parameters_distribution(model):
    num_parameters = dict()
    for name, param in model.named_parameters():
        if param.dim() > 1:
            num_parameters[name] = param.numel()
    fig = plt.figure(figsize=(8, 6))
    plt.grid(axis='y')
    plt.bar(list(num_parameters.keys()), list(num_parameters.values()))
    plt.title('#Parameter Distribution')
    plt.ylabel('Number of Parameters')
    plt.xticks(rotation=60)
    plt.tight_layout()
    plt.show()

plot_num_parameters_distribution(model)
```



Select Sparsity Based on Sensitivity Curves and #Parameters Distribution

Question 5

Based on the sensitivity curves and the distribution of #parameters in the model, please select the sparsity for each layer.

Note that the overall compression ratio of pruned model mostly depends on the layers with larger #parameters, and different layers have different sensitivity to pruning (see Question 4).

Please make sure that after pruning, the sparse model is 25% of the size of the dense model, and validation accuracy is higher than 92.5 after finetuning.

Hint:

- The layer with more #parameters should have larger sparsity. (see *Figure #Parameter Distribution*)
- The layer that is sensitive to the pruning sparsity (i.e., the accuracy will drop quickly as sparsity becomes higher) should have smaller sparsity. (see *Figure Sensitivity Curves*)

```
In [107]: recover_model()

sparsity_dict = {
##### YOUR CODE STARTS HERE #####
    # please modify the sparsity value of each layer
    # please DO NOT modify the key of sparsity_dict
    'backbone.conv0.weight': 0.1,
    'backbone.conv1.weight': 0.1,
    'backbone.conv2.weight': 0.2,
    'backbone.conv3.weight': 0.5,
    'backbone.conv4.weight': 0.6,
    'backbone.conv5.weight': 0.8,
    'backbone.conv6.weight': 0.8,
    'backbone.conv7.weight': 0.9,
    'classifier.weight': 0.1
##### YOUR CODE ENDS HERE #####
}
```

Please run the following cell to prune the model according to your defined `sparsity_dict`, and print the information of sparse model.

```
In [108]: pruner = FineGrainedPruner(model, sparsity_dict)
print(f'After pruning with sparsity dictionary')
for name, sparsity in sparsity_dict.items():
    print(f'    {name}: {sparsity:.2f}')
print(f'The sparsity of each layer becomes')
for name, param in model.named_parameters():
    if name in sparsity_dict:
        print(f'    {name}: {get_sparsity(param):.2f}')

sparse_model_size = get_model_size(model, count_nonzero_only=True)
print(f"Sparse model has size={sparse_model_size / MiB:.2f} MiB = {sparse_model_size / 1024:.2f}% of dense model size")
sparse_model_accuracy = evaluate(model, dataloader['test'])
print(f"Sparse model has accuracy={sparse_model_accuracy:.2f}% before fintuning")

#plot_weight_distribution(model, count_nonzero_only=True)
```

After pruning with sparsity dictionary

```
backbone.conv0.weight: 0.10
backbone.conv1.weight: 0.10
backbone.conv2.weight: 0.20
backbone.conv3.weight: 0.50
backbone.conv4.weight: 0.60
backbone.conv5.weight: 0.80
backbone.conv6.weight: 0.80
backbone.conv7.weight: 0.90
classifier.weight: 0.10
```

The sparsity of each layer becomes

```
backbone.conv0.weight: 0.10
backbone.conv1.weight: 0.10
backbone.conv2.weight: 0.20
backbone.conv3.weight: 0.50
backbone.conv4.weight: 0.60
backbone.conv5.weight: 0.80
backbone.conv6.weight: 0.80
backbone.conv7.weight: 0.90
classifier.weight: 0.10
```

Sparse model has size=8.62 MiB = 24.49% of dense model size

```
eval: 0% | 0/20 [00:00<?, ?it/s]
```

Sparse model has accuracy=88.65% before fintuning

Finetune the fine-grained pruned model

As we can see from the outputs of previous cell, even though fine-grained pruning reduces the most of model weights, the accuracy of model also dropped. Therefore, we have to finetune the sparse model to recover the accuracy.

Please run the following cell to finetune the sparse model. It should take around 3 minutes to finish.

```
In [109]: num_finetune_epochs = 5
optimizer = torch.optim.SGD(model.parameters(), lr=0.01, momentum=0.9, weight_decay=
scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, num_finetune_epochs)
criterion = nn.CrossEntropyLoss()

best_sparse_model_checkpoint = dict()
best_accuracy = 0
print(f'Finetuning Fine-grained Pruned Sparse Model')
for epoch in range(num_finetune_epochs):
    # At the end of each train iteration, we have to apply the pruning mask
    # to keep the model sparse during the training
    train(model, dataloader['train'], criterion, optimizer, scheduler,
          callbacks=[lambda: pruner.apply(model)])
    accuracy = evaluate(model, dataloader['test'])
    is_best = accuracy > best_accuracy
    if is_best:
        best_sparse_model_checkpoint['state_dict'] = copy.deepcopy(model.state_dict())
        best_accuracy = accuracy
    print(f'    Epoch {epoch+1} Accuracy {accuracy:.2f}% / Best Accuracy: {best_accuracy:.2f}%')
```

Finetuning Fine-grained Pruned Sparse Model

train: 0%| | 0/98 [00:00<?, ?it/s]

eval: 0%| | 0/20 [00:00<?, ?it/s]

Epoch 1 Accuracy 92.62% / Best Accuracy: 92.62%

train: 0%| | 0/98 [00:00<?, ?it/s]

eval: 0%| | 0/20 [00:00<?, ?it/s]

Epoch 2 Accuracy 92.64% / Best Accuracy: 92.64%

train: 0%| | 0/98 [00:00<?, ?it/s]

eval: 0%| | 0/20 [00:00<?, ?it/s]

Epoch 3 Accuracy 92.69% / Best Accuracy: 92.69%

train: 0%| | 0/98 [00:00<?, ?it/s]

eval: 0%| | 0/20 [00:00<?, ?it/s]

Epoch 4 Accuracy 92.67% / Best Accuracy: 92.69%

train: 0%| | 0/98 [00:00<?, ?it/s]

eval: 0%| | 0/20 [00:00<?, ?it/s]

Epoch 5 Accuracy 92.83% / Best Accuracy: 92.83%

Run the following cell to see the information of best finetuned sparse model.

```
In [110]: # load the best sparse model checkpoint to evaluate the final performance
model.load_state_dict(best_sparse_model_checkpoint['state_dict'])
sparse_model_size = get_model_size(model, count_nonzero_only=True)
print(f"Sparse model has size={sparse_model_size / MiB:.2f} MiB = {sparse_model_size / dense_model_size:.2f}% of dense model size")
sparse_model_accuracy = evaluate(model, dataloader['test'])
print(f"Sparse model has accuracy={sparse_model_accuracy:.2f}% after fintuning")
```

Sparse model has size=8.62 MiB = 24.49% of dense model size

eval: 0% | 0/20 [00:00<?, ?it/s]

Sparse model has accuracy=92.83% after fintuning

Channel Pruning (OPTIONAL)

In this section, we will implement the channel pruning. Channel pruning removes an entire channel, so that it can achieve inference speed up on existing hardware like GPUs. Similarly, we remove the channels whose weights are of smaller magnitudes (measured by Frobenius norm).

```
In [111]: # firstly, let's restore the model weights to the original dense version
# and check the validation accuracy
recover_model()
dense_model_accuracy = evaluate(model, dataloader['test'])
print(f"dense model has accuracy={dense_model_accuracy:.2f}%")
```

eval: 0% | 0/20 [00:00<?, ?it/s]

dense model has accuracy=92.95%

Remove Channel Weights (OPTIONAL)

Unlike fine-grained pruning, we can remove the weights entirely from the tensor in channel pruning. That is to say, the number of output channels is reduced:

$$\#out_channels_{new} = \#out_channels_{origin} \cdot (1 - sparsity)$$

The weight tensor W is still dense after channel pruning. Thus, we will refer to *sparsity* as **prune ratio**.

Like fine-grained pruning, we can use different pruning rates for different layers. However, we use a uniform pruning rate for all the layers for now. We are targeting 2x computation reduction, which is roughly 30% uniform pruning rate (think about why).

Feel free to try out different pruning ratios per layer at the end of this section. You can pass in a list of ratios to the `channel_prune` function.

Question 6 (OPTIONAL)

Please complete the following functions for channel pruning.

Here we naively prune all output channels other than the first $\#out_channels_{new}$ channels.


```

In [124]: def get_num_channels_to_keep(channels: int, prune_ratio: float) -> int:
    """A function to calculate the number of layers to PRESERVE after pruning
    Note that preserve_rate = 1. - prune_ratio
    """

    ##### YOUR CODE STARTS HERE #####
    return int(round(channels * (1 - prune_ratio)))
    ##### YOUR CODE ENDS HERE #####

@torch.no_grad()
def channel_prune(model: nn.Module,
                  prune_ratio: Union[List, float]) -> nn.Module:
    """Apply channel pruning to each of the conv layer in the backbone
    Note that for prune_ratio, we can either provide a floating-point number,
    indicating that we use a uniform pruning rate for all layers, or a list of
    numbers to indicate per-layer pruning rate.
    """

    # sanity check of provided prune_ratio
    assert isinstance(prune_ratio, (float, list))
    n_conv = len([m for m in model.backbone if isinstance(m, nn.Conv2d)])
    # note that for the ratios, it affects the previous conv output and next
    # conv input, i.e., conv0 - ratio0 - conv1 - ratio1-...
    if isinstance(prune_ratio, list):
        assert len(prune_ratio) == n_conv - 1
    else: # convert float to list
        prune_ratio = [prune_ratio] * (n_conv - 1)

    # we prune the convs in the backbone with a uniform ratio
    model = copy.deepcopy(model) # prevent overwrite
    # we only apply pruning to the backbone features
    all_convs = [m for m in model.backbone if isinstance(m, nn.Conv2d)]
    all_bns = [m for m in model.backbone if isinstance(m, nn.BatchNorm2d)]
    # apply pruning. we naively keep the first k channels
    assert len(all_convs) == len(all_bns)
    for i_ratio, p_ratio in enumerate(prune_ratio):
        prev_conv = all_convs[i_ratio]
        prev_bn = all_bns[i_ratio]
        next_conv = all_convs[i_ratio + 1]
        original_channels = prev_conv.out_channels # same as next_conv.in_channels
        n_keep = get_num_channels_to_keep(original_channels, p_ratio)

        # prune the output of the previous conv and bn
        #print("before pruning prev conv weight shape: ", prev_conv.weight.detach().shape)
        prev_conv.weight.set_(prev_conv.weight.detach()[:n_keep])
        #print("after pruning prev conv weight shape: ", prev_conv.weight.detach().shape)
        prev_bn.weight.set_(prev_bn.weight.detach()[:n_keep])
        prev_bn.bias.set_(prev_bn.bias.detach()[:n_keep])
        prev_bn.running_mean.set_(prev_bn.running_mean.detach()[:n_keep])
        prev_bn.running_var.set_(prev_bn.running_var.detach()[:n_keep])

        # prune the input of the next conv (hint: just one line of code)
        ##### YOUR CODE STARTS HERE #####
        #print("next conv weight shape: ", next_conv.weight.detach().shape)
        next_conv.weight.set_(next_conv.weight.detach()[:, :n_keep])
        #print("after pruning next conv weight shape: ", next_conv.weight.detach().shape)
        ##### YOUR CODE ENDS HERE #####

    return model

```

Run the following cell to perform a sanity check to make sure the implementation is correct.

```
In [125]: dummy_input = torch.randn(1, 3, 32, 32).cuda()
pruned_model = channel_prune(model, prune_ratio=0.3)
pruned_macs = get_model_macs(pruned_model, dummy_input)
assert pruned_macs == 305388064
print('* Check passed. Right MACs for the pruned model.')
```

* Check passed. Right MACs for the pruned model.

Now let's evaluate the performance of the model after uniform channel pruning with 30% pruning rate.

As you may see, directly removing 30% of the channels leads to low accuracy.

```
In [126]: pruned_model_accuracy = evaluate(pruned_model, dataloader['test'])
print(f"pruned model has accuracy={pruned_model_accuracy:.2f}%")
```

eval: 0% | 0/20 [00:00<?, ?it/s]

pruned model has accuracy=28.15%

Ranking Channels by Importance (OPTIONAL)

As you can see, removing the first 30% of channels in all layers leads to significant accuracy reduction. One potential method to remedy the issue is to find the **less important** channel weights to remove. A popular criterion for importance is to use the Frobenius norm of the weights corresponding to each input channel:

$$importance_i = \|W_i\|_2, \quad i = 0, 1, 2, \dots, \text{\#in_channels} - 1$$

We can sort the channel weights from more important to less important, and then keep the first k channels for each layer.

Question 7 (OPTIONAL)

Please complete the following functions for sorting the weight tensor based on the Frobenius norm.

Hint:

- To calculate Frobenius norm of a tensor, Pytorch provides `torch.norm` (<https://pytorch.org/docs/master/generated/torch.norm.html?highlight=torch+norm#torch.norm>) APIs.

```

In [129]: # function to sort the channels from important to non-important
def get_input_channel_importance(weight):
    in_channels = weight.shape[1]
    importances = []
    # compute the importance for each input channel
    for i_c in range(weight.shape[1]):
        channel_weight = weight.detach()[ :, i_c]
        ##### YOUR CODE STARTS HERE #####
        importance = torch.norm(channel_weight)
        ##### YOUR CODE ENDS HERE #####
        importances.append(importance.view(1))
    return torch.cat(importances)

@torch.no_grad()
def apply_channel_sorting(model):
    model = copy.deepcopy(model) # do not modify the original model
    # fetch all the conv and bn layers from the backbone
    all_convs = [m for m in model.backbone if isinstance(m, nn.Conv2d)]
    all_bns = [m for m in model.backbone if isinstance(m, nn.BatchNorm2d)]
    # iterate through conv layers
    for i_conv in range(len(all_convs) - 1):
        # each channel sorting index, we need to apply it to:
        # - the output dimension of the previous conv
        # - the previous BN layer
        # - the input dimension of the next conv (we compute importance here)
        prev_conv = all_convs[i_conv]
        prev_bn = all_bns[i_conv]
        next_conv = all_convs[i_conv + 1]
        # note that we always compute the importance according to input channels
        importance = get_input_channel_importance(next_conv.weight)
        # sorting from large to small
        sort_idx = torch.argsort(importance, descending=True)

        # apply to previous conv and its following bn
        prev_conv.weight.copy_(torch.index_select(
            prev_conv.weight.detach(), 0, sort_idx))
        for tensor_name in ['weight', 'bias', 'running_mean', 'running_var']:
            tensor_to_apply = getattr(prev_bn, tensor_name)
            tensor_to_apply.copy_(
                torch.index_select(tensor_to_apply.detach(), 0, sort_idx)
            )

        # apply to the next conv input (hint: one line of code)
        ##### YOUR CODE STARTS HERE #####
        next_conv.weight.copy_(
            torch.index_select(next_conv.weight.detach(), 1, sort_idx)
        )
        ##### YOUR CODE ENDS HERE #####

    return model

```

Now run the following cell to sanity check if the results are correct.

```
In [130]: print('Before sorting...')
dense_model_accuracy = evaluate(model, dataloader['test'])
print(f"dense model has accuracy={dense_model_accuracy:.2f}%")

print('After sorting...')
sorted_model = apply_channel_sorting(model)
sorted_model_accuracy = evaluate(sorted_model, dataloader['test'])
print(f"sorted model has accuracy={sorted_model_accuracy:.2f}%")

# make sure accuracy does not change after sorting, since it is
# equivalent transform
assert abs(sorted_model_accuracy - dense_model_accuracy) < 0.1
print('* Check passed.')
```

Before sorting...

eval: 0%| | 0/20 [00:00<?, ?it/s]

dense model has accuracy=92.95%

After sorting...

eval: 0%| | 0/20 [00:00<?, ?it/s]

sorted model has accuracy=92.95%

* Check passed.

Finally, we compare the pruned models' accuracy with and without sorting.

```
In [131]: channel_pruning_ratio = 0.3 # pruned-out ratio

print(" * Without sorting...")
pruned_model = channel_prune(model, channel_pruning_ratio)
pruned_model_accuracy = evaluate(pruned_model, dataloader['test'])
print(f"pruned model has accuracy={pruned_model_accuracy:.2f}%")

print(" * With sorting...")
sorted_model = apply_channel_sorting(model)
pruned_model = channel_prune(sorted_model, channel_pruning_ratio)
pruned_model_accuracy = evaluate(pruned_model, dataloader['test'])
print(f"pruned model has accuracy={pruned_model_accuracy:.2f}%")
```

* Without sorting...

eval: 0%| | 0/20 [00:00<?, ?it/s]

pruned model has accuracy=28.15%

* With sorting...

eval: 0%| | 0/20 [00:00<?, ?it/s]

pruned model has accuracy=36.81%

As you can see, the channel sorting can slightly improve the pruned model's accuracy, but there is still a huge degrade, which is quite common for channel pruning. But luckily, we can perform fine-tuning to recover the accuracy.

```
In [132]: num_finetune_epochs = 5
optimizer = torch.optim.SGD(pruned_model.parameters(), lr=0.01, momentum=0.9, weight_decay=0.01)
scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, num_finetune_epochs)
criterion = nn.CrossEntropyLoss()

best_accuracy = 0
for epoch in range(num_finetune_epochs):
    train(pruned_model, dataloader['train'], criterion, optimizer, scheduler)
    accuracy = evaluate(pruned_model, dataloader['test'])
    is_best = accuracy > best_accuracy
    if is_best:
        best_accuracy = accuracy
    print(f'Epoch {epoch+1} Accuracy {accuracy:.2f}% / Best Accuracy: {best_accuracy:.2f}%')
```

```
train:  0%|          | 0/98 [00:00<?, ?it/s]
eval:   0%|          | 0/20 [00:00<?, ?it/s]
Epoch 1 Accuracy 91.62% / Best Accuracy: 91.62%
train:  0%|          | 0/98 [00:00<?, ?it/s]
eval:   0%|          | 0/20 [00:00<?, ?it/s]
Epoch 2 Accuracy 91.68% / Best Accuracy: 91.68%
train:  0%|          | 0/98 [00:00<?, ?it/s]
eval:   0%|          | 0/20 [00:00<?, ?it/s]
Epoch 3 Accuracy 92.02% / Best Accuracy: 92.02%
train:  0%|          | 0/98 [00:00<?, ?it/s]
eval:   0%|          | 0/20 [00:00<?, ?it/s]
Epoch 4 Accuracy 92.22% / Best Accuracy: 92.22%
train:  0%|          | 0/98 [00:00<?, ?it/s]
eval:   0%|          | 0/20 [00:00<?, ?it/s]
Epoch 5 Accuracy 92.33% / Best Accuracy: 92.33%
```

Measure acceleration from pruning (OPTIONAL)

After fine-tuning, the model almost recovers the accuracy. You may have already learned that channel pruning is usually more difficult to recover accuracy compared to fine-grained pruning. However, it directly leads to a smaller model size and smaller computation without specialized model format. It can also run faster on GPUs. Now we compare the model size, computation, and latency of the pruned model.

```

In [133]: # helper functions to measure latency of a regular PyTorch models.
# Unlike fine-grained pruning, channel pruning
# can directly leads to model size reduction and speed up.
@torch.no_grad()
def measure_latency(model, dummy_input, n_warmup=20, n_test=100):
    model.eval()
    # warmup
    for _ in range(n_warmup):
        _ = model(dummy_input)
    # real test
    t1 = time.time()
    for _ in range(n_test):
        _ = model(dummy_input)
    t2 = time.time()
    return (t2 - t1) / n_test # average latency

table_template = "{:<15} {:<15} {:<15} {:<15}"
print (table_template.format('', 'Original', 'Pruned', 'Reduction Ratio'))

# 1. measure the latency of the original model and the pruned model on CPU
# which simulates inference on an edge device
dummy_input = torch.randn(1, 3, 32, 32).to('cpu')
pruned_model = pruned_model.to('cpu')
model = model.to('cpu')

pruned_latency = measure_latency(pruned_model, dummy_input)
original_latency = measure_latency(model, dummy_input)
print(table_template.format('Latency (ms)',
                             round(original_latency * 1000, 1),
                             round(pruned_latency * 1000, 1),
                             round(original_latency / pruned_latency, 1)))

# 2. measure the computation (MACs)
original_macs = get_model_macs(model, dummy_input)
pruned_macs = get_model_macs(pruned_model, dummy_input)
print(table_template.format('MACs (M)',
                             round(original_macs / 1e6),
                             round(pruned_macs / 1e6),
                             round(original_macs / pruned_macs, 1)))

# 3. measure the model size (params)
original_param = get_num_parameters(model)
pruned_param = get_num_parameters(pruned_model)
print(table_template.format('Param (M)',
                             round(original_param / 1e6, 2),
                             round(pruned_param / 1e6, 2),
                             round(original_param / pruned_param, 1)))

# put model back to cuda
pruned_model = pruned_model.to('cuda')
model = model.to('cuda')

```

	Original	Pruned	Reduction Ratio
Latency (ms)	7.0	4.9	1.4
MACs (M)	606	305	2.0
Param (M)	9.23	5.01	1.8

Question 8 (OPTIONAL)

Please answer the following questions using the information in the previous code cell.

Question 8.1 (OPTIONAL)

Explain why removing 30% of channels roughly leads to 50% computation reduction.

Your Answer:

Question 8.2 (OPTIONAL)

Explain why the latency reduction ratio is slightly smaller than computation reduction.

Your Answer:

Compare Fine-grained Pruning and Channel Pruning (OPTIONAL)

Question 9 (OPTIONAL)

After all experiments in this lab, you may have become familiar with both fine-grained pruning and channel pruning.

Please answer the following questions using what you have learned from the lectures and this lab.

Question 9.1 (OPTIONAL)

What are the advantages and disadvantages of fine-grained pruning and channel pruning? You can discuss from the perspective of compression ratio, accuracy, latency, hardware support (*i.e.*, requiring specialized hardware accelerator), etc.

Your Answer:

Question 9.2 (OPTIONAL)

If you want to make your model run faster on a smartphone, which pruning method will you use? Why?

Your Answer:

