

11-695: AI Engineering Devices Management II

LTI/SCS

Spring 2020

① Py-torch for tensorflow users

② Data Paralellism with pytorch

③ `torch.distributed`

torch_general_model

```
1 from torch import nn
2 import torch.nn.functional as F
3
4 class Mnist_CNN(nn.Module):
5     def __init__(self):
6         super().__init__()
7         self.conv1 = nn.Conv2d(1, 16, kernel_size=3, stride=2, padding=1)
8         self.conv2 = nn.Conv2d(16, 16, kernel_size=3, stride=2, padding=1)
9         self.conv3 = nn.Conv2d(16, 10, kernel_size=3, stride=2, padding=1)
10
11     def forward(self, xb):
12         xb = xb.view(-1, 1, 28, 28)
13         xb = F.relu(self.conv1(xb))
14         xb = F.relu(self.conv2(xb))
15         xb = F.relu(self.conv3(xb))
16         xb = F.avg_pool2d(xb, 4)
17         return xb.view(-1, xb.size(1))
```

- Like `tf.keras`, we can either write the whole network as a class extending from `nn.Module`

¹https://pytorch.org/tutorials/beginner/nn_tutorial.html

torch_general_model_2

```
1 class Lambda(nn.Module):
2     def __init__(self, func):
3         super().__init__()
4         self.func = func
5     def forward(self, x):
6         return self.func(x)
7
8 model = nn.Sequential(
9     Lambda(lambda x: x.view(-1, 1, 28, 28)),
10    nn.Conv2d(1, 16, kernel_size=3, stride=2, padding=1),
11    nn.ReLU(),
12    nn.Conv2d(16, 16, kernel_size=3, stride=2, padding=1),
13    nn.ReLU(),
14    nn.Conv2d(16, 10, kernel_size=3, stride=2, padding=1),
15    nn.ReLU(),
16    nn.AvgPool2d(4),
17    Lambda(lambda x: x.view(x.size(0), -1)),
18 )
```

- Or can use `nn.Sequential`

torch_autograd

```
1 model = torch.nn.Sequential(  
2     torch.nn.Linear(D_in, H),  
3     torch.nn.ReLU(),  
4     torch.nn.Linear(H, D_out))  
5 loss_fn = torch.nn.MSELoss(reduction='sum')  
6  
7 for t in range(epochs):  
8     x, y = from_dataset(...)  
9     y_pred = model(x)  
10    loss = loss_fn(y_pred, y)  
11  
12    # backprop calculation  
13    model.zero_grad() # manually zero out all Grads  
14    loss.backward() # new Grads calculation  
15  
16    # SGD update  
17    with torch.no_grad(): # meaning not tracking grads change  
18        for param in model.parameters():  
19            param -= learning_rate * param.grad
```

- They have a same purpose: auto-tracking grads

²https://pytorch.org/tutorials/beginner/pytorch_with_examples.html

torch.optim

```
1 model = torch.nn.Sequential(  
2     torch.nn.Linear(D_in, H),  
3     torch.nn.ReLU(),  
4     torch.nn.Linear(H, D_out))  
5 loss_fn = torch.nn.MSELoss(reduction='sum')  
6 opt = torch.optim.Adam(model.parameters(), lr=1e03)  
7  
8 for t in range(epochs):  
9     x, y = from_dataset(...)  
10    y_pred = model(x)  
11    loss = loss_fn(y_pred, y)  
12  
13    # backprop calculation  
14    optimizer.zero_grad() # Note: not model.zero_grad()  
15    loss.backward() # same  
16  
17    # SGD update  
18    opt.step() # simple
```

- Use `torch.optim` for convenience

tf_devices

```
1 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
2
3 # move model to device
4 net.to(device)
5
6 # move every input batch to device
7 inputs, labels = data[0].to(device), data[1].to(device)
8
9 # optimize as usual
10 loss = criterion(model(inputs), labels)
11 ... # backprop and SGD (more later)
12
13 # transfer back to CPU
14 loss.cpu().numpy() # Note: for weights w it usually has w.grad and w.data
15
16 # or along with recasting dtype
17 loss.to("cpu", torch.double))
```

- Remember to reside your model in GPUs first to accept inputs

³https://pytorch.org/tutorials/beginner/deep_learning_60min_blitz.html

torch_dynamic

```
1 class DynamicNet(torch.nn.Module):
2     def __init__(self, D_in, H, D_out):
3         super(DynamicNet, self).__init__()
4         self.input_linear = torch.nn.Linear(D_in, H)
5         self.middle_linear = torch.nn.Linear(H, H)
6         self.output_linear = torch.nn.Linear(H, D_out)
7
8     def forward(self, x):
9         """Reuse the middle_linear Module random times
10
11         Here we also see that it is perfectly safe to reuse the same Module many
12         times when defining a computational graph. This is a big improvement from Lua
13         Torch, where each Module could be used only once."""
14         h_relu = self.input_linear(x).clamp(min=0)
15         for _ in range(random.randint(0, 3)):
16             h_relu = self.middle_linear(h_relu).clamp(min=0)
17         y_pred = self.output_linear(h_relu)
18         return y_pred
```

- This dynamic design is one of the biggest advantages on Pytorch

⁴https://pytorch.org/tutorials/beginner/pytorch_with_examples.html

torch_dataset

```
1 from torch.utils.data import DataLoader
2
3 # you have tensors already
4 train_ds = TensorDataset(x_train, y_train)
5 train_dl = DataLoader(train_ds, batch_size=32, shuffle=True, num_workers=4)
6
7 # or from a raw data folder
8 data_transform = transforms.Compose([
9     transforms.RandomSizedCrop(224),
10    transforms.RandomHorizontalFlip(),
11    transforms.ToTensor(),
12    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
13 ])
14 my_dataset = datasets.ImageFolder(root='my_data/train',
15                                 transform=data_transform)
16 dataset_loader = torch.utils.data.DataLoader(my_dataset,
17                                             batch_size=4, shuffle=True,
18                                             num_workers=4)
```

- Similar to `tf.keras`

⁵https://pytorch.org/tutorials/beginner/data_loading_tutorial.html

torch_train_eval

```
1 criterion = nn.CrossEntropyLoss()
2 opt = optim.SGD(model.parameters(), lr=0.001, momentum=0.9)
3 for epoch in range(epochs):
4     model.train()
5     for xb, yb in train_dl:
6         pred = model(xb)
7         loss = criterion(pred, yb)
8
9         # almost always do this triplet:
10        opt.zero_grad() # have to do this MANUALLY
11        loss.backward() # backprop: calculate grads
12        opt.step() # one SGD-based update step
13
14
15    model.eval() # turn on this flag for evaluation/test
16    with torch.no_grad(): # don't forget to turn off grads tracking as well
17        valid_loss = sum(loss_func(model(xb), yb) for xb, yb in valid_dl)
18
19    print(epoch, valid_loss / len(valid_dl))
```

- For training: `model.train()`
- Otherwise: `model.eval()`, with `torch.no_grad()`

① Py-torch for tensorflow users

② Data Paralellism with pytorch

③ `torch.distributed`

torch_data_parallel_multiple_gpus

```
1 # feed data into a data loader as normal
2 train_ds = TensorDataset(x_train, y_train)
3 train_dl = DataLoader(train_ds, batch_size=32, shuffle=True, num_workers=4)
4
5 # define model as normal
6 model = MNIST_CNN() # see slide no. 3
7
8 # signify data parallelism
9 if torch.cuda.device_count() > 1:
10     model = nn.DataParallel(model)
11
12 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
13 model.to(device)
14
15 for data in rand_loader:
16     input = data.to(device) # 4 nodes will have 8 data samples each
17     output = model(input)
```

- pytorch splits data equally for us
- It will collect and merge results before returning

⁶https://pytorch.org/tutorials/beginner/blitz/data_parallel_tutorial.html

- ① Py-torch for tensorflow users
- ② Data Paralellism with pytorch
- ③ `torch.distributed`

python multiprocessing

```
1 from multiprocessing import Process
2
3 def f(pid):
4     print('hello ', pid)
5
6 if __name__ == '__main__':
7     ps = []
8     for i in range(5):
9         p = Process(target=f, args=(i,))
10        p.start()
11        ps.append(p)
12
13    for p in ps:
14        p.join() # wait until completion
```

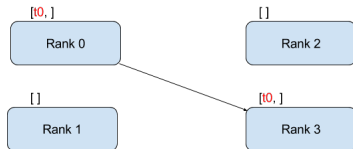
- Usually each process will process a segment of a big job
- In advanced settings, processes can communicate with each other
- Restricted to one machine

⁷<https://docs.python.org/3.6/library/multiprocessing.html>

torch_dist_multiprocessing

```
1 import torch.distributed as dist
2 from torch.multiprocessing import Process
3
4 def run(rank, size): # to be passed into a single process later
5     pass
6
7 def init_process(rank, size, fn, backend='gloo'):
8     """ Initialize the distributed environment. """
9     os.environ['MASTER_ADDR'] = '127.0.0.1'
10    os.environ['MASTER_PORT'] = '29500'
11    dist.init_process_group(backend, rank=rank, world_size=size)
12    fn(rank, size)
13
14 if __name__ == "__main__":
15     processes = []
16     for rank in range(5):
17         p = Process(target=init_process, args=(rank, size, run))
18         p.start()
19         processes.append(p)
20     for p in processes:
21         p.join()
```

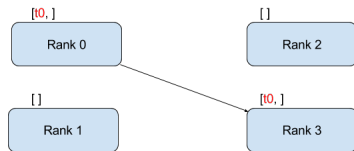
- torch.distributed extends into multiple machines/clusters



torch_dist_send_recv_sync

```
1 def run(rank, size):
2     tensor = torch.zeros(1)
3     if rank == 0:
4         tensor += 1
5         dist.send(tensor=tensor, dst=1) # Send the tensor to process 1
6     else:
7         dist.recv(tensor=tensor, src=0) # Receive tensor from process 0
8     print('Rank ', rank, ' has data ', tensor[0])
```

- Processes wait until communication ends to resume

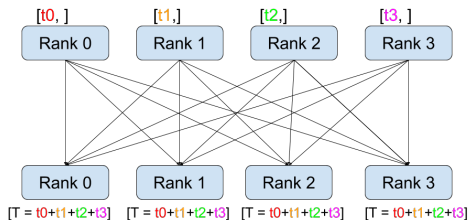


`torch_dist_send_recv_Async`

```
1 def run(rank, size):
2     tensor = torch.zeros(1)
3     req = None
4     if rank == 0:
5         tensor += 1
6         req = dist.isend(tensor=tensor, dst=1) # Send the tensor to process 1
7     else:
8         req = dist.irecv(tensor=tensor, src=0) # Receive tensor from process 0
9     req.wait()
10    print('Rank ', rank, ' has data ', tensor[0])
```

- Processes does not have to wait, might cause inconsistencies

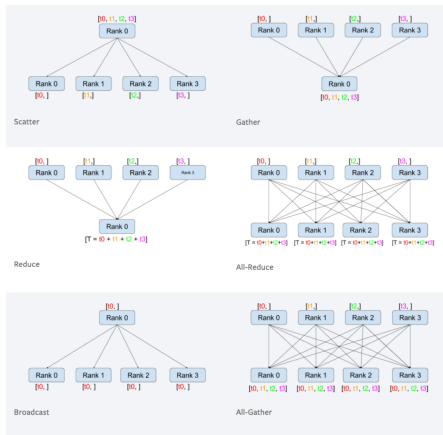
Image credit: Séb Arnold



`torch_dist_allreduce`

```
1 """ Gradient averaging. """
2 def average_gradients(model):
3     size = float(dist.get_world_size())
4     for param in model.parameters():
5         dist.all_reduce(param.grad.data, op=dist.reduce_op.SUM)
6         param.grad.data /= size
```

- All to all with reduce ops: SUM, PRODUCT, MAX, MIN



- Also can also use a checkpoint/barrier: `dist.barrier()`

torch_dist_SGD

```
1 def run(rank, size):
2     torch.manual_seed(1234)
3     train_set, bsz = partition_dataset()
4     model = Net()
5     optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.5)
6
7     num_batches = ceil(len(train_set.dataset) / float(bsz))
8     for epoch in range(10):
9         epoch_loss = 0.0
10        for data, target in train_set:
11            optimizer.zero_grad()
12            output = model(data)
13            loss = F.nll_loss(output, target)
14            epoch_loss += loss.item()
15            loss.backward()
16            average_gradients(model)
17            optimizer.step()
18        print('Rank ', dist.get_rank(), ', epoch ', epoch, ': ', epoch_loss / num_batches)
```

- `torch.dist` helps distribute data to backend nodes/ranks

⁸https://pytorch.org/tutorials/intermediate/dist_tuto.html

Backend	gloo		mpi		nccl	
	CPU	GPU	CPU	GPU	CPU	GPU
send	✓	✗	✓	?	✗	✗
recv	✓	✗	✓	?	✗	✗
broadcast	✓	✓	✓	?	✗	✓
all_reduce	✓	✓	✓	?	✗	✓
reduce	✓	✗	✓	?	✗	✓
all_gather	✓	✗	✓	?	✗	✓
gather	✓	✗	✓	?	✗	✗
scatter	✓	✗	✓	?	✗	✗
barrier	✓	✗	✓	?	✗	✓

- gloo: CPU and some (not optimized for) GPUs (Facebook)
- nccl: optimized for CUDA/GPU (NVIDIA)
- mpi: support popular platforms: Open-MPI, MVAPICH2, Intel MPI. Needs extra work to get it installed.

⁹<https://pytorch.org/docs/stable/distributed.html>