## **Network Training**

**Baotong Tian** 

ECE 208/408 – The Art of Machine Learning

(Slides adapted from Huiran Yu's version last year, and <u>http://cs231n.stanford.edu/slides/2022/lecture 7 ruohan.pdf</u> as well as official documentation from PyTorch and NVIDIA. )

## Popular Deep Learning Frameworks



#### **Computational graph**

Imperative: Imperative-style programs perform computation as you run them





Gluon: new MXNet interface to accelerate research



NumP	y
------	---

import numpy as np
np.random.seed(0)

N, D = 3, 4

x = np.random.randn(N, D) y = np.random.randn(N, D) z = np.random.randn(N, D) a = x \* y

b = a + z c = np.sum(b) grad\_c = 1.0 grad\_b = grad\_c \* np.ones((N, D)) grad\_a = grad\_b.copy() grad\_z = grad\_b.copy() grad\_x = grad\_a \* y grad\_y = grad\_a \* x PyTorch

import torch

print(x.grad)

N, D = 3, 4
x = torch.randn(N, D, requires\_grad=True)
y = torch.randn(N, D, requires\_grad=True)
z = torch.randn(N, D, requires\_grad=True)
a = x \* y
b = a + z
c = torch.sum(b)
c.backward()
print(x.grad)
print(x.grad)

X

#### Computational graph

Ζ

Define **tensor** with gradient required, which will be added to the computational graph

#### PyTorch

#### import torch

N, D = 3, 4

x	=	torch.randn(N,	D,	requires_grad=True)
У	=	<pre>torch.randn(N,</pre>	D,	requires_grad=True)
z	=	torch.randn(N,	D,	requires_grad=True)

```
a = x * y
b = a + z
c = torch.sum(b)
c.backward()
print(x.grad)
print(x.grad)
print(x.grad)
```

#### Computational graph



The forward pass looks just like numpy.

Remember, the function \*, +, torch.sum() here are **pytorch functions**, these functions will build the dependency between tensors.

#### PyTorch

#### import torch

```
N, D = 3, 4
x = torch.randn(N, D, requires_grad=True)
y = torch.randn(N, D, requires_grad=True)
z = torch.randn(N, D, requires_grad=True)
```

```
a = x * y

b = a + z

c = torch.sum(b)
```

```
c.backward()
```

print(x.grad)
print(x.grad)
print(x.grad)

#### Computational graph



Recall HW5, we do the backpropagation manually, but pytorch can do it automatically, due to the computational graphs.

c.backward() will calculate the derivative of c with respect to x, y, z, and will write the gradient to the grad attribute of x, y, z.

#### PyTorch

import torch

N, D = 3, 4
x = torch.randn(N, D, requires\_grad=True)
y = torch.randn(N, D, requires\_grad=True)

z = torch.randn(N, D, requires\_grad=True)

a = x \* y b = a + zc = torch.sum(b)

c.backward()	
print(x.grad)	
print(x.grad)	
print(x.grad)	

#### Module

A neural network layer is a module, such as a convolution layer (torch.nn.Conv2d). It often has the following:

- Weight attribute: store the weight of convolution kernel.
- Weight Initialization method: initialize the weight.
- **Forward** method: the computation build in the forward part
- **Backward** : this part is invisible to users, but is implemented by Pytorch already.
- **Buffer** attribute: weights but don't require grad

## Pytorch: Two levels of abstraction

Tensor:

- if Tensor.requires\_grad==False, it is imperative ndarray, but no gradient will be computed

- if **Tensor.requires\_grad==True**, it is a node in a computational graph; stores data and gradient

Module:

neural network layer(s); store learnable weights.

#### Module

#### ⊟ torch.nn

Parameters

#### □ Convolution Layers

Conv1d Conv2d Conv3d ConvTranspose1d

ConvTranspose2d

ConvTranspose3d

Other layers:

Dropout, Linear,

Normalization Layer

#### □ torch.nn Parameters

Containers
Convolution Layers

#### Pooling Layers

MaxPool1d MaxPool2d MaxPool3d MaxUnpool1d MaxUnpool2d MaxUnpool3d AvgPool1d AvgPool2d AvgPool3d FractionalMaxPool2d LPPool2d AdaptiveMaxPool1d AdaptiveMaxPool2d AdaptiveMaxPool3d AdaptiveAvgPool1d AdaptiveAvgPool2d

#### □ Loss functions L1Loss **MSELoss** CrossEntropyLoss NLLLoss PoissonNLLLoss **KLDivLoss** BCELoss BCEWithLogitsLoss MarginRankingLoss HingeEmbeddingLoss MultiLabelMarginLoss SmoothL1Loss SoftMarginLoss MultiLabelSoftMarginLoss CosineEmbeddingLoss MultiMarginLoss **TripletMarginLoss**

Network Training, ECE 208/408 - The Art of Machine Learning, Spring 2025

AdaptiveAvgPool3d



## Module

class Net(nn.Module):

```
def __init__(self):
    super(Net, self).__init__()
    self.conv1 = nn.Conv2d(1, 10, kernel_size=5)
    self.conv2 = nn.Conv2d(10, 20, kernel_size=5)
    self.mp = nn.MaxPool2d(2)
    self.fc = nn.Linear(320, 10) # 320 -> 10
    def forward(self, x):
        in_size = x.size(0)
        x = F.relu(self.mp(self.conv1(x)))
        x = F.relu(self.mp(self.conv2(x)))
        x = x.view(in_size, -1) # fLatten the tensor
        x = self.fc(x)
        return F.log_softmax(x)
Network Training, ECE 208/408 - The Art of Machine Learning, Spring 2025
```

## Define a CNN

- Net class is a CNN defined by user, it inherits from torch.nn.Module
- Initialize the basic layers in \_\_init\_\_. Usually we only use the basic layer provided by pytorch to build our own network
- Define the forward method to build your computational graph
- When you call a module, it will automatically call the forward method of the module.

#### Check the model structure

#### torchinfo · PyPI

from torchinfo import summary

model = Net()
batch\_size = 16
summary(model, input\_size=(batch\_size, 1, 28, 28))

Layer (type:depth-idx)	Output Shape	Param #
	[16, 10]	
├─Conv2d: 1-1	[16, 10, 24, 24]	260
-MaxPool2d: 1-2	[16, 10, 12, 12]	
-Conv2d: 1-3	[16, 20, 8, 8]	5,020
-MaxPool2d: 1-4	[16, 20, 4, 4]	
Linear: 1-5	[16, 10]	3,210
Total params: 8,490 Trainable params: 8,490 Non-trainable params: 0 Total mult-adds (M): 7.59		
Input size (MB): 0.05 Forward/backward pass size (MB): 0.90 Params size (MB): 0.03 Estimated Total Size (MB): 0.99		

```
Dataset and Dataloader
```

```
class CustomImageDataset(Dataset):
   def __init__(self, annotations_file, img_dir, transform=None, target_transform=None):
        self.img_labels = pd.read_csv(annotations_file)
        self.img_dir = img_dir
        self.transform = transform
        self.target_transform = target_transform
   def len (self):
       return len(self.img_labels)
   def __getitem__(self, idx):
        img_path = os.path.join(self.img_dir, self.img_labels.iloc[idx, 0])
        image = read_image(img_path)
       label = self.img_labels.iloc[idx, 1]
        if self.transform:
            image = self.transform(image)
        if self.target_transform:
           label = self.target_transform(label)
        sample = {"image": image, "label": label}
       return sample
```

#### **Dataset and Dataloader**

The Dataset retrieves our dataset's features and labels one sample at a time. While training a model, we typically want to pass samples in "minibatches", reshuffle the data at every epoch to reduce model overfitting, and use Python's multiprocessing to speed up data retrieval.

DataLoader is an iterable that abstracts this complexity for us in an easy API.

```
from torch.utils.data import DataLoader
```

```
train_dataloader = DataLoader(training_data, batch_size=64, shuffle=True)
test_dataloader = DataLoader(test_data, batch_size=64, shuffle=True)
```

P

## Define a loss function and optimizer

Put net.parameter() to optim.SGD, so the gradient descent can be applied to the parameters of CNN.

CrossEntropyLoss() is also a module.

import torch.optim as optim
criterion = nn.CrossEntropyLoss()
optimizer = optim.SGD(net.parameters(), lr=0.001, momentum=0.9)

Adam optimizer is recommended. Reference to GBC Ch. 8.5

```
Train a neural network
```

set the model to the training mode, but it will not do any training. It inform layers such as Dropout and BatchNorm

```
https://stackoverflow.com/questions
/51433378/what-does-model-train-d
o-in-pytorch
```

```
def train(dataloader, model, loss_fn, optimizer):
                      size = len(dataloader.dataset)
                     model.train()
                      for batch, (X, y) in enumerate(dataloader):
                          X, y = X.to(device), y.to(device)
                          # Compute prediction error
                          pred = model(X)
                          loss = loss_fn(pred, y)
                          # Backpropagation
                          optimizer.zero grad()
                          loss.backward()
                          optimizer.step()
                          if batch % 100 == 0:
                              loss, current = loss.item(), (batch + 1) \star len(X)
                              print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")
Network Training, ECE 208/408 - The Art of Machine Learning, Spring 2025
                                                                                       17
```

Iterate the dataloader to get mini-batches of data

```
def train(dataloader, model, loss_fn, optimizer):
    size = len(dataloader.dataset)
    model.train()
    for batch, (X, y) in enumerate(dataloader):
        X, y = X.to(device), y.to(device)
```

```
# Compute prediction error
pred = model(X)
loss = loss_fn(pred, y)
# Backpropagation
optimizer.zero_grad()
loss.backward()
optimizer.step()
if batch % 100 == 0:
loss, current = loss.item(), (batch + 1) * len(X)
print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")
Network Training, ECE 208/408 - The Art of Machine Learning, Spring 2025
```

Equivalent to calling model.forward(X)

Compute loss

```
def train(dataloader, model, loss_fn, optimizer):
    size = len(dataloader.dataset)
    model.train()
    for batch, (X, y) in enumerate(dataloader):
        X, y = X.to(device), y.to(device)
```

```
# Compute prediction error
pred = model(X)
loss = loss_fn(pred, y)
```

```
# Backpropagation
optimizer.zero_grad()
loss.backward()
optimizer.step()
```

```
if batch % 100 == 0:
    loss, current = loss.item(), (batch + 1) * len(X)
    print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")
```

#### Train a neural network def train(dataloader, model, loss\_fn, optimizer): size = len(dataloader.dataset) model.train() for batch, (X, y) in enumerate(dataloader): X, y = X.to(device), y.to(device) *# Compute prediction error* pred = model(X)loss = loss\_fn(pred, y) # Backpropagation set the gradients to zero optimizer.zero grad() loss.backward() optimizer.step() https://stackoverflow.com/questions if batch % 100 == 0: /48001598/why-do-we-need-to-callloss, current = loss.item(), (batch + 1) \* len(X) zero-grad-in-pytorch print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")

```
def train(dataloader, model, loss_fn, optimizer):
    size = len(dataloader.dataset)
    model.train()
    for batch, (X, y) in enumerate(dataloader):
        X, y = X.to(device), y.to(device)
```

```
# Compute prediction error
pred = model(X)
loss = loss_fn(pred, y)
```

loss.backward() calculate all the gradient of loss w.r.t parameters # Backpropagation

optimizer.zero\_grad()
loss.backward()

optimizer.step()

```
if batch % 100 == 0:
    loss, current = loss.item(), (batch + 1) * len(X)
    print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")
```

```
def train(dataloader, model, loss_fn, optimizer):
    size = len(dataloader.dataset)
    model.train()
    for batch, (X, y) in enumerate(dataloader):
        X, y = X.to(device), y.to(device)
        # Compute prediction error
        pred = model(X)
        loss = loss_fn(pred, y)
        # Backpropagation
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```

optimizer.step() will update the parameter using SGD: weight = weight - lr\*weight.grad

```
if batch % 100 == 0:
    loss, current = loss.item(), (batch + 1) * len(X)
    print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")
```

loss.item() is to convert a

```
def train(dataloader, model, loss_fn, optimizer):
                                       size = len(dataloader.dataset)
                                       model.train()
                                       for batch, (X, y) in enumerate(dataloader):
                                           X, y = X.to(device), y.to(device)
                                           # Compute prediction error
                                           pred = model(X)
                                           loss = loss_fn(pred, y)
                                           # Backpropagation
                                           optimizer.zero_grad()
                                           loss.backward()
                                           optimizer.step()
                                           if batch % 100 == 0:
                                               loss, current = loss.item(), (batch + 1) * len(X)
torch.float type scaler into float
                                               print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")
```

## Save and load the trained network

When the training reaches the end or some particular conditions, for example the highest precision in the validation set, we would like to save the current parameters of the model. For more information, please check: <u>Saving and Loading</u> <u>Models - PyTorch Tutorials</u> <u>save:</u>

torch.save(model.state\_dict(), PATH)

Load:

model = TheModelClass(\*args, \*\*kwargs)
model.load\_state\_dict(torch.load(PATH))
model.eval()

## Save and Load the trained network: Takeaways

- When loading models, use map\_location='cpu' in most cases to conserve VRAM.
- In distributed training, ensure that checkpoints are saved only on the master process.
- To resume training, include the optimizer state in the checkpoint. For inference purposes, omit the optimizer state to reduce storage requirements.

## PyTorch tutorials

https://pytorch.org/tutorials/beginner/basics/quickstart\_tutorial.html

http://cs231n.stanford.edu/slides/2022/discussion\_4\_pytorch.pdf

https://pytorch.org/tutorials/beginner/blitz/cifar10\_tutorial.html

#### **GPU** acceleration



https://transfer.d2.mpi-inf.mpg.de/rs hetty/hlcv/Pytorch\_tutorial.pdf

## Using GPU

Bluehive: BluehiveInfo.pdf

Google Colab: Runtime / change runtime type

device = 'cuda' if torch.cuda.is\_available() else 'cpu'
Move Tensors and Modules to GPU: .cuda() or .to(device)
Monitor GPU usage: nvidia-smi or gpustat · PyPI

# How to prevent **overfitting** when training deep neural networks?

#### Three perspectives

Data

Model

Training strategies

## Data augmentation

Augment the training data



https://www.baeldung.com/cs/ml-data-augmentation

#### Data preprocessing

Make the optimization more stable and easier



#### Reduce model complexity

Reduce the number of layers or neurons in the network

Use simpler activation functions.



#### Add regularization term to the loss function

$$\widehat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} \quad \underbrace{J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y})}_{(i)} + \underbrace{\lambda}_{(iii)} \underbrace{R(\boldsymbol{\theta})}_{(ii)}.$$

(i) fit the training data

(ii) the regularization term, such as L1 or L2 norm.

(iii) hyperparameter for controlling the trade-off



Tools for observing learning curves: tensorboard, wandb

#### Tensorboard

- Install tensorboard with pip or conda
- Add code in your training script:



• Run "tensorboard --logdir path/to/logdir" with command line and it will provide you a link to access the tensorboard web application. You can also launch it in the jupyter notebook if you need.

#### Tensorboard



## **Batch normalization**

During training, it normalizes the activation values across the batch.

During testing, it uses the mean and variance values determined in the training.



## Dropout

During training, in each forward pass, randomly set some neurons to zero. Probability of dropping is a hyperparameter; 0.5 is common



At test time, all probability.

## How to choose hyperparameters when training deep neural networks?

## Choose hyperparameters

Choose the batch size according to your device.

Start with a learning rate that makes training loss go down. If not, overfit a small batch of samples to debug.

Look at the learning curves (loss and metrics).

## Learning rate



LWLS Figure 5-7

#### Learning curves



#### Training Acceleration: torch.compile



Training Acceleration: torch.compile

```
def foo(x, y):
    a = torch.sin(x)
    b = torch.cos(y)
    return a + b
opt_foo1 = torch.compile(foo)
print(opt_foo1(torch.randn(10, 10), torch.randn(10, 10)))
```

## Training Acceleration: Mixed Precision TRAINING

#### MIXED PRECISION TRAINING



#### **Training Acceleration: Mixed Precision**

```
# Creates model and optimizer in default precision
model = Net().cuda()
optimizer = optim.SGD(model.parameters(), ...)
# Creates a GradScaler once at the beginning of training.
scaler = GradScaler()
for epoch in epochs:
    for input, target in data:
        optimizer.zero grad()
        # Runs the forward pass with autocasting.
        with autocast(device_type='cuda', dtype=torch.float16):
            output = model(input)
            loss = loss fn(output, target)
        # Scales loss. Calls backward() on scaled loss to create scaled gradients.
        # Backward passes under autocast are not recommended.
        # Backward ops run in the same dtype autocast chose for corresponding forward ops.
        scaler.scale(loss).backward()
        # scaler.step() first unscales the gradients of the optimizer's assigned params.
        # If these gradients do not contain infs or NaNs, optimizer.step() is then called,
        # otherwise, optimizer.step() is skipped.
        scaler.step(optimizer)
        # Updates the scale for next iteration.
        scaler.update()
```

Network Training, ECE 208/408 - The Art of Machine Learning, Spring 2025

ſ

#### Lecture wrap up

We covered **PyTorch** basics. Practice with homework 6.

Preventing **overfitting** in deep neural networks requires a combination of techniques, including using more data, regularization, early stopping, reducing model complexity, dropout, batch normalization, etc.

Use learning curves to tune hyperparameters.