# Lec17: Neural Networks

#### Introduction

- (Deep) Neural Network is a really flexible and powerful machine learning model that
  - Approximate any reasonable function in theory
  - "Learns the pattern from data" incredibly well in practice, particularly in hard areas like cognition, perception
- Two computational modes:
  - Training: go through training examples, update model parameters (usually driven by reducing a certain loss function)
  - Inference: given new input, compute output (prediction, next token, classification, regressions, probability, etc...)
- Usually very hungry for data+compute
  - Modern LLM pretty much demands to be run on accelerators (GPUs, TPUs, ...)

### Agenda

- To study how a neural network works, in particular its computational structure and their acceleration on GPU...
- We study a relatively simple but representative example: digit classifier on the MNIST dataset.
- Moving from higher abstraction -> lower abstraction
  - PyTorch version
  - Python+numpy version (pytorch free) + Math
  - C version (without dependencies)
  - C+CUDA version (accelerated on GPU!)

#### Neural Network Background

• 10,000-ft bird view:

A particular class of mathematical function

 $f(x; \mathbf{w}) = y,$ 

x is input vector w is the model parameters y is output vector

#### • Training:

given (x0,y0), (x1,y1), ..., update w to make f() better predict y

 Inferencing: given x, compute f(x; w) = y

## Function f(x;w)

- It's constructed in layers
  input x -> (layer 1) -> output z1 -> (layer 2) -> output z2 -> ...
  -> output z
- There are many different layers:
  - Linear layer:  $f(x) := x \cdot W + b = z$ W is weights, a matrix of size (#input, #output), b is vector bias.
  - Convolutional layer, recurrent layer, embedding layer...
  - Activation layer: ReLU:  $f(x) = \max(x, 0)$ Sigmoid, Softmax, ...
  - Normalization layer (to stabilize training): BatchNorm, LayerNorm,
  - Attention layer:

# **MNIST** Dataset

- MNIST contains 60,000 training images and 10,000 testing images of handwritten digits.
- The dataset comprises grayscale images of size 28×28 pixels.

000000000000000000  $\mathbf{X}$ 2222222222222222 3333333333333333333 448444444444 Ч *777777777777* 8888888888888888 8 B в **99999999999** 9 999

## Simple Multilayer Perceptron (MLP)

- Input x: 784 (28\*28) pixels(floats)
- Three layers: input -> hidden -> output
  - Linear layer1: 784\*256 (a matrix that maps 784-vector to 256-vector z1=x\*W1+b1
  - ReLU layer: 256->256 a1=ReLU(z1)
  - Linear layer2: 256\*10 (a matrix that maps 256 –vector to 10-vector) z2=a1W2+b2

**Softmax + Cross-Entropy Loss:** 

 $L = -\sum y \log(\operatorname{softmax}(z2))$ 



| ♥ | | ♥.7 \ <

#### PyTorch MLP: Network



#### PyTorch MLP: Train

```
# Training the model
def train(model, criterion, optimizer, epoch):
    model.train()
    running_loss = 0.0
    for i in range(iters_per_epoch):
        optimizer.zero_grad()
        data = train_data[i * batch_size : (i + 1) * batch_size].to("cuda")
        target = train_labels[i * batch_size : (i + 1) * batch_size].to("cuda")
        start = time.time()
        outputs = model(data)
        loss = criterion(outputs, target)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
        end = time.time()
        running_loss += loss.item()
        if i % 100 == 99 or i == 0:
            print(f"Epoch: {epoch+1}, Iter: {i+1}, Loss: {loss}")
            print(f"Iteration Time: {(end - start) * 1e3:.4f} sec")
            running loss = 0.0
```

#### PyTorch: Evaluate

```
# Evaluation function to report average batch accuracy using the loaded test data
def evaluate(model, test_data, test_labels):
    device = torch.device("cuda" if torch.cuda.is available() else "cpu")
    model.to(device)
    model.eval()
    total_batch_accuracy = torch.tensor(0.0, device=device)
    num_batches = 0
    with torch.no grad():
        for i in range(len(test_data) // batch_size):
            data = test_data[i * batch_size : (i + 1) * batch_size].to(device)
            target = test_labels[i * batch_size : (i + 1) * batch_size].to(device)
            outputs = model(data)
            _, predicted = torch.max(outputs.data, 1)
            correct batch = (predicted == target).sum().item()
            total batch = target.size(0)
            if total batch != 0: # Check to avoid division by zero
                batch accuracy = correct batch / total batch
                total_batch_accuracy += batch_accuracy
                num batches += 1
    avg_batch_accuracy = total_batch_accuracy / num_batches
    print(f"Average Batch Accuracy: {avg batch accuracy * 100:.2f}%")
```

### Python MLP without PyTorch

- What actually does PyTorch do?
  - Layers are just functions mapping input vector -> output vector
  - .backward(): automatically computing gradients using chain rules
  - Train: iteration of computing gradient and update parameters using some variant of stochastic gradient descent
- For our simple MLP we can derive and hardcode the gradient computing rule, the forward pass of the network, and the training iterations with SGD using nothing other than Numpy linear algebra

### **Training: Back-propagation**

- Recall training: f(x; w) = yOptimize w to minimize (crossentropy) loss function over a training data set (x0, y0), (x1, y1), ..., (xn, yn):  $\min_{w} \sum_{i}^{i} L(f(x_{i}; w), y_{i})$
- Take one training point (x0, y0). How to update w so that loss is lower?
- Intuition: In calculus, we learn that gradient (derivative) at w is the slope that the function increases around w

def L(w): return 3\*w\*\*2 - 4\*w + 5



#### **Gradient & Back-propagation**

• Intuitively, the f gradient of w at point x0  $\nabla_w f(x0, w)$ 

points to the **direction** in which the scalar function f(x0,w) increases the fastest;

• If we want to **reduce** f(x0,w), we would go tiny step in opposite direction of the gradient  $g := \nabla_w f(x0,w)$ :

$$\delta w = t * g$$

- $f(x0; w \delta w) \approx f(x0; w) (\delta w, g) = f(x0; w) t ||g||^2$
- So in summary, to minimize the loss function at point (x0,y0), one would update parameters w a tiny amount in the  $\nabla L_w(x0; w)$
- But how to compute  $\nabla L_w(x0; w)$  for our MLP? Chain rule, and layer by layer backwards (that's where back-propagation got its name)

#### Chain rule

**General Chain Rule Principle** 

For a composite function:

L = f(g(h(x)))

The gradient of L w.r.t. x is:

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial h} \cdot \frac{\partial h}{\partial x}$$

#### Chain rule applied to the 3-layer net

**Forward Pass:** 

$Z_1 = XW_1 + b_1$	(Layer 1)
$A_1 = \operatorname{ReLU}(Z_1)$	(Activation)
$Z_2 = A_1 W_2 + b_2$	(Layer 2)
$A_2 = \operatorname{Softmax}(Z_2)$	(Output)
$L=\mathrm{CrossEntropy}(A_2,y)$	(Loss)

**Backward Pass (Gradient Flow):** 

1. Output Layer Gradient:

$$rac{\partial L}{\partial Z_2} = A_2 - y \quad ext{(Softmax + CrossEntropy derivative)}$$

2. Layer 2 Weights & Bias:

$$rac{\partial L}{\partial W_2} = A_1^T rac{\partial L}{\partial Z_2}, \quad rac{\partial L}{\partial b_2} = \sum rac{\partial L}{\partial Z_2}$$

3. Gradient to Layer 1 Output:

$$rac{\partial L}{\partial A_1} = rac{\partial L}{\partial Z_2} W_2^T$$

4. Layer 1 Activation Gradient:

$$rac{\partial L}{\partial Z_1} = rac{\partial L}{\partial A_1} \odot \mathrm{ReLU}'(Z_1)$$

5. Layer 1 Weights & Bias:

$$rac{\partial L}{\partial W_1} = X^T rac{\partial L}{\partial Z_1}, \quad rac{\partial L}{\partial b_1} = \sum rac{\partial L}{\partial Z_1}$$

**General Chain Rule Principle** 

For a composite function:

$$L = f(g(h(x)))$$

The gradient of L w.r.t. x is:

 $rac{\partial L}{\partial x} = rac{\partial L}{\partial f} \cdot rac{\partial f}{\partial g} \cdot rac{\partial g}{\partial h} \cdot rac{\partial h}{\partial x}$ 

### Layer-by-layer Back-propagation

• For an linear layer,

z = xW + b

the backward() should compute the following gradients: given gradient\_out  $(\partial L/\partial z)$ , computes gradients for parameters W, b, and for gradient\_in( $\partial L/\partial x$ ):

For a linear layer Z = XW + b, the gradients are:

• 
$$\frac{\partial L}{\partial W} = X^T \frac{\partial L}{\partial Z}$$
  
•  $\frac{\partial L}{\partial b} = \sum \frac{\partial L}{\partial Z}$  (summed over the batch)  
•  $\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Z} W^T$ 

## Numpy version

<u>https://github.com/Infatoshi/mnist-cuda/blob/master/python/c-friendly.py</u>

#### **References:**

- <u>https://github.com/Infatoshi/mnist-cuda/tree/master</u>
- <u>https://www.kaggle.com/code/scaomath/simple-neural-network-</u> <u>for-mnist-numpy-from-scratch/notebook</u>
- Back-propagation: Andrej Karpathy, MicroGrad: https://www.youtube.com/watch?v=VMj-3S1tku0&t=579s
- Pure CUDA implementation to reproduce GPT-2(124M) model: <u>https://github.com/karpathy/llm.c/discussions/481</u> <u>https://www.youtube.com/watch?v=l8pRSuU81PU</u>