

Neural Networks and Deep Learning

Lecture Notes

Based on Andrej Karpathy's Lectures

Andrej Karpathy
OpenAI / Tesla

August 19, 2025

These notes are based on Andrej Karpathy's lectures on neural networks and deep learning fundamentals.

“The goal is to build a strong foundation in neural networks by implementing everything from scratch and understanding the math.”

Contents

1	Lecture 1: The Spelled-Out Intro to Neural Networks and Backpropagation with Micrograd	3
1.1	Introduction to Neural Network Training and Micrograd	3
1.1.1	Micrograd: A Pedagogical Autograd Engine	3
1.2	Understanding Derivatives: The Intuition	3
1.2.1	Derivative of a Single-Variable Function	3
1.2.2	Derivative with Multiple Inputs (Partial Derivatives)	5
1.3	Building the Value Object in Micrograd	5
1.3.1	Basic Value Class Structure	6
1.3.2	Overloading Operators for Graph Construction	6
1.3.3	Visualizing the Expression Graph	9
1.4	Backpropagation: The Automatic Gradient Algorithm	9
1.4.1	The Chain Rule	9
1.4.2	How Backpropagation Works in Practice	9
1.5	Building a Neural Network (MLP)	10
1.5.1	The Neuron	10
1.5.2	The Layer	11
1.5.3	The Multi-Layer Perceptron (MLP)	11
1.6	Neural Network Training Loop (Gradient Descent)	12
1.6.1	The Loss Function	12
1.6.2	The Training Loop Steps	13
1.6.3	Learning Rate and Stability	14
1.7	Micrograd vs. PyTorch: A Comparison	14
1.7.1	Similarities in API and Core Concepts	14
1.7.2	Key Differences and Production-Grade Features	15
1.8	Conclusion	16

1 Lecture 1: The Spelled-Out Intro to Neural Networks and Backpropagation with Micrograd

Abstract

This lecture provides a thorough deep dive into the foundational concepts of neural network training, focusing on automatic differentiation (autograd) and backpropagation. Using a simplified Python library called Micrograd, we will build a neural network from scratch, demystifying the "under the hood" mechanisms. The pedagogical approach emphasizes understanding scalar operations and the chain rule, avoiding the complexities of high-dimensional tensors initially to promote intuitive grasp. We will cover derivatives, the Value object, manual and automatic backpropagation, the training loop, and compare Micrograd's approach to production-grade libraries like PyTorch.

1.1 Introduction to Neural Network Training and Micrograd

Neural network training is fundamentally about iteratively tuning the weights of a neural network to minimize a loss function, thereby improving the network's accuracy. This process relies heavily on an algorithm called backpropagation, which efficiently calculates the gradient of the loss function with respect to the network's weights.

1.1.1 Micrograd: A Pedagogical Autograd Engine

Micrograd is a Python library designed to illustrate the core principles of automatic gradient computation (autograd) and backpropagation. Its primary goal is pedagogical:

- It operates on scalar values (single numbers) rather than complex N-dimensional tensors commonly found in modern deep learning libraries like PyTorch or JAX. This simplification allows for a clearer understanding of the underlying mathematical operations and the chain rule without being bogged down by tensor dimensionality.
- It is intentionally concise, with the core autograd engine (responsible for backpropagation) implemented in roughly 100 lines of very simple Python code. The neural network library built on top of it (nn.py) is also quite minimal, defining basic structures like neurons, layers, and multi-layer perceptrons (MLPs).

While Micrograd is not for production use due to its scalar-based nature and lack of parallelization, it provides a crucial foundation for understanding how modern deep learning frameworks function at their mathematical core.

1.2 Understanding Derivatives: The Intuition

Before diving into backpropagation, it's essential to have a strong intuitive understanding of what a derivative represents.

1.2.1 Derivative of a Single-Variable Function

A derivative measures the sensitivity of a function's output to a tiny change in its input. It tells us the slope of the function at a specific point, indicating whether the function is increasing or decreasing and by how much.

Consider a scalar-valued function $f(x) = 3x^2 - 4x + 5$. We can numerically approximate the derivative at a point x using the definition:

$$\frac{df}{dx} \approx \frac{f(x+h) - f(x)}{h}$$

where h is a very small number (e.g., 0.001).

```

1  import math
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  # Define the function
6  def f(x):
7      return 3*x**2 - 4*x + 5
8
9  # Example: calculate f(3.0)
10 print(f(3.0)) # Output: 20.0
11
12 # Plotting the function
13 xs = np.arange(-5, 5, 0.25)
14 ys = f(xs)
15 plt.plot(xs, ys)
16 plt.title("Function f(x) = 3x^2 - 4x + 5")
17 plt.xlabel("x")
18 plt.ylabel("f(x)")
19 plt.grid(True)
20 plt.show()
21
22 # Numerical derivative at x = 3.0
23 h = 0.001
24 x = 3.0
25 f_x = f(x)
26 f_x_plus_h = f(x + h)
27 slope_at_3 = (f_x_plus_h - f_x) / h
28 print(f"Slope at x={x}: {slope_at_3}") # Expected: ~14.0 (
    Analytical: 6x - 4 -> 6*3 - 4 = 14)
29
30 # Numerical derivative at x = -3.0
31 x_neg = -3.0
32 f_x_neg = f(x_neg)
33 f_x_plus_h_neg = f(x_neg + h)
34 slope_at_neg_3 = (f_x_plus_h_neg - f_x_neg) / h
35 print(f"Slope at x={x_neg}: {slope_at_neg_3}") # Expected: ~-
    22.0 (Analytical: 6*(-3) - 4 = -22)

```

Listing 1: Numerical Derivative Calculation

The sign of the derivative indicates the direction of change:

- Positive derivative: Function increases if input is slightly increased.
- Negative derivative: Function decreases if input is slightly increased.
- Zero derivative: Function is momentarily flat (e.g., at a minimum or maximum).

1.2.2 Derivative with Multiple Inputs (Partial Derivatives)

When a function has multiple inputs, we talk about partial derivatives. A partial derivative with respect to one input tells us how the output changes when only that specific input is slightly nudged, while all other inputs are held constant.

Consider the expression $d = a \times b + c$ where $a = 2$, $b = -3$, $c = 10$. The output $d = 4$. Let's find the partial derivative of d with respect to a , i.e., $\frac{\partial d}{\partial a}$:

```

1 # Define the multi-input function
2 def func_d(a, b, c):
3     return a * b + c
4
5 # Initial values
6 a, b, c = 2.0, -3.0, 10.0
7 d1 = func_d(a, b, c) # d1 = 4.0
8
9 h = 0.001
10
11 # d(d)/d(a)
12 a_bumped = a + h
13 d2_a = func_d(a_bumped, b, c)
14 slope_a = (d2_a - d1) / h
15 print(f"d(d)/d(a): {slope_a}") # Expected: -3.0 (which is the
    value of b)
16
17 # d(d)/d(b)
18 b_bumped = b + h
19 d2_b = func_d(a, b_bumped, c)
20 slope_b = (d2_b - d1) / h
21 print(f"d(d)/d(b): {slope_b}") # Expected: 2.0 (which is the
    value of a)
22
23 # d(d)/d(c)
24 c_bumped = c + h
25 d2_c = func_d(a, b, c_bumped)
26 slope_c = (d2_c - d1) / h
27 print(f"d(d)/d(c): {slope_c}") # Expected: 1.0 (coefficient
    of c)

```

Listing 2: Numerical Partial Derivative Calculation

This numerical verification matches the analytical derivatives:

- $\frac{\partial d}{\partial a} = b = -3$
- $\frac{\partial d}{\partial b} = a = 2$
- $\frac{\partial d}{\partial c} = 1 = 1$

This intuition is critical: the derivative tells us the direct influence of a small change in an input on the output.

1.3 Building the Value Object in Micrograd

Micrograd's core data structure is the Value object, which wraps a scalar number and tracks how it was computed, forming an expression graph.

1.3.1 Basic Value Class Structure

The Value class holds the actual numerical data and a grad attribute, which will store the derivative of the final output (loss) with respect to this value. Initially, grad is set to zero, implying no influence on the output until gradients are computed.

```

1 class Value:
2     def __init__(self, data, _children=(), _op=''):
3         self.data = data
4         self.grad = 0.0 # Initialize gradient to zero
5         # Internal variables to build the expression graph
6         self._prev = set(_children) # Set of Value objects
7         # that are children
8         self._op = _op # String representing the operation
9         # that produced this value
10
11     def __repr__(self):
12         # Provides a nice string representation for printing
13         return f"Value(data={self.data}, grad={self.grad})"
14
15 # Example usage
16 a = Value(2.0)
17 print(a) # Output: Value(data=2.0, grad=0.0)

```

Listing 3: Initial Value Class

1.3.2 Overloading Operators for Graph Construction

To build mathematical expressions naturally (e.g., $a + b$, $a * b$), we overload Python's special methods (dunder methods) like `__add__` and `__mul__`. Crucially, these methods not only perform the operation but also store the "lineage" of the new Value object: its `_prev` (children nodes) and `_op` (operation name).

```

1 import math
2
3 class Value:
4     def __init__(self, data, _children=(), _op=''):
5         self.data = data
6         self.grad = 0.0
7         self._prev = set(_children)
8         self._op = _op
9         # Stores the local backward function for this node
10        self._backward = lambda: None # Default empty
11        # backward function for leaf nodes
12
13    def __repr__(self):
14        return f"Value(data={self.data}, grad={self.grad})"
15
16    def __add__(self, other):
17        # Handle scalar addition (e.g., Value + 1)
18        other = other if isinstance(other, Value) else Value(
19            other)
20        out = Value(self.data + other.data, (self, other), '+')

```

```

19
20     def _backward():
21         # For addition, the gradient is simply passed
22           through (local derivative is 1)
23         self.grad += out.grad * 1.0 # Accumulate gradient
24         other.grad += out.grad * 1.0 # Accumulate
25           gradient
26     out._backward = _backward
27     return out
28
29     def __mul__(self, other):
30         # Handle scalar multiplication (e.g., Value * 2)
31         other = other if isinstance(other, Value) else Value(
32           other)
33     out = Value(self.data * other.data, (self, other), '**'
34           )
35
36     def _backward():
37         # For multiplication, local derivatives are the '
38           other' operand
39         self.grad += out.grad * other.data # Accumulate
40           gradient
41         other.grad += out.grad * self.data # Accumulate
42           gradient
43     out._backward = _backward
44     return out
45
46     # For operations like 2 * a (reversed multiplication)
47     def __rmul__(self, other):
48         return self * other
49
50     def __pow__(self, other):
51         # This implementation expects 'other' to be a scalar (
52           int or float), not a Value object
53         assert isinstance(other, (int, float)), "only
54           supporting int/float powers for now"
55         out = Value(self.data**other, (self,), f'**{other}')
56
57         def _backward():
58             # Power rule: d/dx(x^n) = n*x^(n-1)
59             self.grad += out.grad * (other * (self.data**(
60               other-1)))
61         out._backward = _backward
62         return out
63
64     def __neg__(self): # -self
65         return self * -1
66
67     def __sub__(self, other): # self - other
68         return self + (-other)
69
70     def __truediv__(self, other): # self / other
71         # Division is implemented as multiplication by a
72           negative power
73         return self * other**-1

```

```

64     def exp(self):
65         x = self.data
66         out = Value(math.exp(x), (self,), 'exp')
67
68         def _backward():
69             # d/dx(e^x) = e^x
70             self.grad += out.grad * out.data
71         out._backward = _backward
72         return out
73
74     def tanh(self):
75         x = self.data
76         t = (math.exp(2*x) - 1) / (math.exp(2*x) + 1)
77         out = Value(t, (self,), 'tanh')
78
79         def _backward():
80             # d/dx(tanh(x)) = 1 - tanh(x)^2
81             self.grad += out.grad * (1 - t**2)
82         out._backward = _backward
83         return out
84
85     def backward(self):
86         # Zero out all gradients first (crucial for iterative
87         # training)
88         # This is a common bug: forgetting to zero gradients
89         # In a real training loop, this would be handled
90         # globally for all parameters
91         # For this Value object, we only zero its own and its
92         # children's gradients for demonstration.
93         # In the full training loop, all network parameters
94         # would be zeroed.
95
96         topo = []
97         visited = set()
98         def build_topo(v):
99             if v not in visited:
100                 visited.add(v)
101                 for child in v._prev:
102                     build_topo(child)
103             topo.append(v)
104         build_topo(self)
105
106         self.grad = 1.0 # Initialize the gradient of the root
107                          # node to 1.0
108
109         for node in reversed(topo):
110             node._backward() # Call the stored backward
111                              # function for each node

```

Listing 4: Adding Arithmetic Operations to Value

Important Considerations in Value Implementation

- **Accumulation of Gradients ($+=$):** When a Value object is used multiple times in an expression (e.g., $b = a + a$), its gradient should be accumulated rather than

overwritten. This is why `self.grad += ...` is used instead of `self.grad =` This is a common bug if not handled correctly.

- **Handling Scalar Operands:** To allow operations like `a + 1` where `1` is a Python `int` (not a `Value` object), the other operand is wrapped in a `Value` object if it's not already one.
- **Reversed Operations (`_rmul_`):** Python's operator precedence means `2 * a` would try to call `2._mul_(a)`, which fails for an `int`. Implementing `_rmul_` (reversed multiply) in `Value` allows Python to fall back to `a._rmul_(2)`, which then correctly calls `a._mul_(2)`.
- **Arbitrary Function `_backward`:** The `_backward` function stored in each `Value` object captures the local derivative calculation. This allows Micrograd to support any operation (like `tanh` or `exp`) as long as its local derivative is known and implemented.

1.3.3 Visualizing the Expression Graph

A helper function, `drawdot`, (not provided in source for brevity, but demonstrated) uses `Graphviz` to visualize the computational graph built by Micrograd, showing nodes (`Value` objects) and edges (operations). This helps to intuitively understand the flow of computation.

1.4 Backpropagation: The Automatic Gradient Algorithm

Backpropagation is the algorithm that efficiently calculates the gradients by recursively applying the chain rule backwards through the computation graph.

1.4.1 The Chain Rule

The chain rule is the mathematical foundation of backpropagation. It allows us to compute the derivative of a composite function. If a variable Z depends on Y , and Y depends on X , then Z also depends on X through Y . The chain rule states:

$$\frac{dZ}{dX} = \frac{dZ}{dY} \times \frac{dY}{dX}$$

Intuitively, if a car travels twice as fast as a bicycle, and the bicycle is four times as fast as a walking man, then the car travels $2 \times 4 = 8$ times as fast as the man. The "rates of change" multiply.

1.4.2 How Backpropagation Works in Practice

1. **Forward Pass:** The mathematical expression is evaluated from inputs to output, computing the data value for each `Value` node.
2. **Initialization:** The grad of the final output (often the loss function) is set to 1.0 (since $\frac{dL}{dL} = 1$). All other grad values are initialized to zero.
3. **Topological Sort:** The nodes of the expression graph are ordered such that a node's children always appear before it in the list (if traversed forward). For backpropagation, this list is then reversed, ensuring that when `_backward` is called on a node, all its dependencies (nodes "further down" the graph) have already had their grad contributions propagated.

```

1  # Assuming 'self' is the root node (e.g., the loss)
2  topo = []
3  visited = set()
4  def build_topo(v):
5      if v not in visited:
6          visited.add(v)
7          for child in v._prev:
8              build_topo(child)
9          topo.append(v)
10 build_topo(self) # Populates 'topo' list
11
12 # Gradients are then computed by iterating in reverse
   order
13 for node in reversed(topo):
14     node._backward()

```

Listing 5: Topological Sort Logic (within backward method)

4. **Backward Pass (`_backward` calls):** Starting from the output node (gradient of 1.0) and iterating backward through the topologically sorted list, each node's `_backward` function is called. This function applies the chain rule: it takes the current node's grad (which is $\frac{dL}{d(\text{current_node})}$) and multiplies it by the *local derivative* of how its children influenced it. The result is then *added* (`+=`) to the children's grad attributes.

- For a `+` operation (e.g., `c = a + b`): $\frac{\partial c}{\partial a} = 1$, $\frac{\partial c}{\partial b} = 1$. So, `a.grad += c.grad * 1.0`, `b.grad += c.grad * 1.0`.
- For a `*` operation (e.g., `c = a * b`): $\frac{\partial c}{\partial a} = b$, $\frac{\partial c}{\partial b} = a$. So, `a.grad += c.grad * b.data`, `b.grad += c.grad * a.data`.
- For `tanh` (e.g., `o = tanh(n)`): $\frac{\partial o}{\partial n} = 1 - \tanh(n)^2 = 1 - o^2$. So, `n.grad += o.grad * (1 - o.data**2)`.

1.5 Building a Neural Network (MLP)

Neural networks, specifically Multi-Layer Perceptrons (MLPs), are just complex mathematical expressions. We can build them using our Value objects.

1.5.1 The Neuron

A neuron is the fundamental building block. It takes multiple inputs (`x`), multiplies them by corresponding weights (`w`), sums these products, adds a bias (`b`), and passes the result through an activation function (e.g., `tanh`).

$$\text{output} = \text{activation_fn} \left(\sum_i (x_i \times w_i) + b \right)$$

```

1  import random
2
3  class Neuron:
4      def __init__(self, nin):
5          # nin: number of inputs to this neuron

```

```

6         self.w = [Value(random.uniform(-1,1)) for _ in range(
            nin)]
7         self.b = Value(random.uniform(-1,1))
8
9     def __call__(self, x):
10        # x: list of input Values
11        # w * x + b (dot product)
12        act = sum((wi*xi for wi, xi in zip(self.w, x)), self.b
13                )
14        out = act.tanh() # Activation function
15        return out
16
17    def parameters(self):
18        # Returns all learnable parameters (weights and bias)
19        return self.w + [self.b]

```

Listing 6: Neuron Class Implementation

1.5.2 The Layer

A layer in an MLP is simply a collection of independent neurons, all receiving the same inputs from the previous layer or the network's initial input.

```

1 class Layer:
2     def __init__(self, nin, nout):
3         # nin: number of inputs to the layer (from previous
4           layer or network input)
5         # nout: number of neurons in this layer (number of
6           outputs from this layer)
7         self.neurons = [Neuron(nin) for _ in range(nout)]
8
9     def __call__(self, x):
10        # x: list of input Values from the previous layer/
11           input
12        # Evaluates each neuron independently
13        outs = [n(x) for n in self.neurons]
14        # Return a list of outputs, or single output if only
15           one neuron
16        return outs if len(outs) == 1 else outs
17
18    def parameters(self):
19        # Collects all parameters from all neurons in this
20           layer
21        return [p for neuron in self.neurons for p in neuron.
22                parameters()]

```

Listing 7: Layer Class Implementation

1.5.3 The Multi-Layer Perceptron (MLP)

An MLP is a sequence of layers, where the outputs of one layer serve as the inputs to the next.

```

1 class MLP:

```

```

2  def __init__(self, nin, nouts):
3      # nin: number of inputs to the MLP
4      # nouts: list of integers defining the sizes of each
        layer
5      # Example: nouts=[4, 4, 1] creates two hidden layers
        of 4 neurons, and an output layer of 1 neuron
6      sz = [nin] + nouts
7      self.layers = [Layer(sz[i], sz[i+1]) for i in range(
        len(nouts))]
8
9  def __call__(self, x):
10     # x: list of input Values to the MLP
11     # Feeds outputs of one layer as inputs to the next
12     for layer in self.layers:
13         x = layer(x)
14     return x
15
16 def parameters(self):
17     # Collects all parameters from all layers in the MLP
18     return [p for layer in self.layers for p in layer.
        parameters()]

```

Listing 8: MLP Class Implementation

1.6 Neural Network Training Loop (Gradient Descent)

The training process involves repeatedly calculating the network's output, measuring its error, and updating its weights using the gradients obtained via backpropagation.

1.6.1 The Loss Function

The loss function quantifies how "bad" the neural network's predictions are compared to the desired targets (ground truth). The goal of training is to minimize this loss. A common example is Mean Squared Error (MSE) loss:

$$L = \frac{1}{N} \sum_{i=1}^N (y_{\text{pred},i} - y_{\text{true},i})^2$$

where y_{pred} are the network's predictions and y_{true} are the actual targets.

```

1  # Example dataset (4 inputs, 4 targets for binary
        classification)
2  xs = [
3      [2.0, 3.0, -1.0],
4      [3.0, -1.0, 0.5],
5      [0.5, 1.0, 1.0],
6      [1.0, 1.0, -1.0],
7  ]
8  ys = [1.0, -1.0, -1.0, 1.0] # Desired targets
9
10 # Initialize a sample MLP
11 # 3 inputs -> 2 hidden layers of 4 neurons each -> 1 output
        neuron
12 n = MLP(3, [4, 4, 1])

```

```

13
14 # Forward pass to get predictions
15 ypred = [n(x) for x in xs]
16 print("Initial predictions:", [p.data for p in ypred])
17
18 # Calculate Mean Squared Error loss
19 loss = sum([(yout - ygt)**2 for ygt, yout in zip(ys, ypred)])
20 print("Initial Loss:", loss.data)

```

Listing 9: Example Data and Loss Calculation

1.6.2 The Training Loop Steps

The core training loop for gradient descent consists of these iterative steps:

1. **Forward Pass:** Feed the input data through the neural network to obtain predictions (ypred).
2. **Calculate Loss:** Compare ypred with ys (targets) using the chosen loss function to get a single loss value.
3. **Zero Gradients:** Before performing backpropagation for the current step, it is crucial to reset all accumulated gradients in the network's parameters to zero. Forgetting this is a common bug, as gradients from previous steps would otherwise incorrectly accumulate.
4. **Backward Pass:** Call `loss.backward()`. This propagates the gradient of the loss all the way back through the network, filling the `grad` attribute of every `Value` object, especially the parameters (weights `w` and biases `b`).
5. **Update Parameters:** Adjust each parameter's data value by taking a small step in the direction opposite to its gradient. This is done to minimize the loss.

$$p.data \leftarrow p.data - \text{learning_rate} \times p.grad$$

Here, `learning_rate` (or `step_size`) is a small positive scalar that controls the magnitude of the update.

```

1 # Re-initialize the network for a fresh start
2 n = MLP(3, [4, 4, 1])
3
4 # Training hyperparameters
5 num_epochs = 50 # Number of training iterations
6 learning_rate = 0.05
7
8 for k in range(num_epochs):
9     # 1. Forward pass
10    ypred = [n(x) for x in xs]
11
12    # 2. Calculate loss
13    loss = sum([(yout - ygt)**2 for ygt, yout in zip(ys, ypred)])
14
15    # 3. Zero gradients (VERY IMPORTANT!)
16    for p in n.parameters():

```

```
17         p.grad = 0.0
18
19     # 4. Backward pass
20     loss.backward()
21
22     # 5. Update parameters (gradient descent step)
23     for p in n.parameters():
24         p.data += -learning_rate * p.grad
25
26     # Print progress
27     print(f"Epoch {k}: Loss = {loss.data:.4f}")
28
29 # Final predictions after training
30 print("\nFinal predictions:", [p.data for p in ypred])
31 print("Desired targets:", ys)
```

Listing 10: Full Training Loop Example

1.6.3 Learning Rate and Stability

The learning_rate (or step_size) is a critical hyperparameter.

- If too low, training will be very slow and may take too long to converge.
- If too high, the optimization can become unstable, oscillate, or even "explode" the loss because it oversteps the optimal direction indicated by the local gradient.

Finding the right learning rate is often an art, though more advanced optimizers automate some of this tuning.

1.7 Micrograd vs. PyTorch: A Comparison

Micrograd's design directly mirrors the core functionalities of modern deep learning frameworks like PyTorch, allowing for a deep intuitive understanding before dealing with production complexities.

1.7.1 Similarities in API and Core Concepts

- **Value vs. torch.Tensor:** Both wrap numerical data and contain a .grad attribute to store gradients. PyTorch's Tensor is an N-dimensional array of scalars, while Micrograd's Value is strictly scalar-valued.
- **Graph Construction:** Both frameworks build a computation graph implicitly as operations are performed on their respective data structures (Value or Tensor).
- **.backward() Method:** Both provide a .backward() method on the root node (e.g., loss) to trigger the backpropagation algorithm.
- **_backward / local_derivative:** The principle of defining how to backpropagate through each atomic operation (by providing its local derivative) is fundamental to both. In PyTorch, custom functions require implementing both a forward and backward method.
- **zero_grad():** The necessity to reset gradients before a new backward pass is a core concept in both Micrograd (manually implemented) and PyTorch (e.g., via optimizer.zero_grad() or model.zero_grad()).

```
1 import torch
2
3 # Similar inputs and weights as in Micrograd neuron example
4 x1 = torch.tensor(2.0, requires_grad=True)
5 w1 = torch.tensor(-3.0, requires_grad=True)
6 x2 = torch.tensor(0.0, requires_grad=True)
7 w2 = torch.tensor(1.0, requires_grad=True)
8 b = torch.tensor(6.8813735870195432, requires_grad=True)
9
10 # Graph construction (forward pass)
11 n = x1*w1 + x2*w2 + b
12 o = n.tanh()
13
14 print(f"PyTorch Forward Pass: {o.item():.7f}")
15
16 # Backward pass
17 o.backward()
18
19 # Gradients
20 print(f"PyTorch Gradients:")
21 print(f" x1.grad: {x1.grad.item():.7f}") # Corresponds to x1.
    grad in Micrograd
22 print(f" w1.grad: {w1.grad.item():.7f}") # Corresponds to w1.
    grad in Micrograd
23 print(f" x2.grad: {x2.grad.item():.7f}") # Corresponds to x2.
    grad in Micrograd
24 print(f" w2.grad: {w2.grad.item():.7f}") # Corresponds to w2.
    grad in Micrograd
25 print(f" b.grad: {b.grad.item():.7f}")   # Corresponds to b.
    grad in Micrograd
26
27 # PyTorch output will match Micrograd's (0.7071067, -
    1.5000000, 1.0000000, 0.0000000, 0.5000000, 0.5000000)
```

Listing 11: PyTorch Equivalent Example for a Neuron

1.7.2 Key Differences and Production-Grade Features

- **Tensors and Efficiency:** PyTorch's primary advantage is its use of tensors, which enable highly optimized, parallelized operations on GPUs, making computations vastly faster than scalar operations for large datasets.
- **Graph Complexity:** Production libraries manage extremely large and dynamic computation graphs, whereas Micrograd's graph construction and traversal are simpler due to its scalar nature.
- **Advanced Optimizers:** While Micrograd uses basic gradient descent, PyTorch offers a wide array of sophisticated optimizers (e.g., Adam, RMSprop) that adapt the learning rate during training.
- **Loss Functions:** PyTorch includes many specialized loss functions (e.g., Cross-Entropy Loss for classification, Max Margin Loss), often with built-in numerical stability features.
- **Batching and Regularization:** PyTorch supports processing data in batches

(subsets of the full dataset) for efficiency with large datasets and includes features like L2 regularization to prevent overfitting.

- **Learning Rate Schedules:** Advanced techniques like learning rate decay, where the learning rate decreases over time, are common in PyTorch to stabilize training towards the end.
- **Implementation Details:** Finding the exact backward pass code for specific operations in a production library like PyTorch can be challenging due to the sheer size and complexity of the codebase, which is highly optimized for performance across different hardware (CPU/GPU kernels) and data types.

1.8 Conclusion

This lecture has provided a comprehensive understanding of neural network training "under the hood" through the lens of Micrograd. We've seen that:

- Neural networks are complex mathematical expressions.
- Training involves minimizing a loss function through iterative gradient descent.
- Backpropagation, a recursive application of the chain rule, efficiently computes these gradients.
- The core Value object and its overloaded operators build the computational graph.
- Understanding the intuitive meaning of derivatives and the chain rule is paramount.

Micrograd, despite its simplicity, demonstrates the fundamental principles that power even the largest neural networks with billions or trillions of parameters used in complex applications like GPT models. The core concepts of forward pass, backward pass (gradient calculation), and parameter updates remain identical, regardless of scale.