

Deep Learning Architectures for Scientific Machine Learning

A Comprehensive PyTorch Guide

Professor Vectorex Gradiens

Wednesday 3rd December, 2025

Version 1.0

Disclaimer

Important Notice: Apart from this paragraph, every single word in this document was written with the assistance of Large Language Models. This document is intended as an aid for learning PyTorch, in particular for implementing deep learning architectures presented in the course "AI in the Sciences and Engineering" held at ETH Zürich.

The content is extremely likely to contain errors, so use it with care, knowing that this isn't reviewed or approved by domain experts. Always verify critical implementations and consult primary sources for production use.

Giovanni Guidarini, December 2025

License: This work is shared for educational purposes. PyTorch is developed by Meta AI and is licensed under the BSD-style license. All code examples are provided as-is without warranty.

Acknowledgments: This textbook builds upon the collective wisdom of the PyTorch community, scientific machine learning researchers, and educators worldwide. Special thanks to the developers of PyTorch, the authors of foundational papers in deep learning, and students whose questions shaped this material.

Preface

My dear student,

Welcome. It is a distinct pleasure to guide you on this journey. You are already embedded in the world of scientific computing, a field of precision, rigor, and computational elegance. You will find that deep learning is not so different. At its heart, it is a powerful form of high-dimensional function approximation, built on the familiar foundations of linear algebra and calculus. The "magic" is simply that we have found a way to make these function approximators (*neural networks*) trainable on a massive scale.

Your lack of PyTorch experience is not a hindrance; it is an opportunity. You arrive with no bad habits. We will build your knowledge from the ground up, with the same care one takes in formulating a proof or designing a robust simulation. PyTorch is our "language" for this. It is expressive, powerful, and, once you are fluent, a genuine joy to use.

This textbook is comprehensive, spanning from the fundamental `torch.Tensor` to state-of-the-art architectures like Transformers and Neural Operators. Each chapter builds systematically, combining rigorous theory with practical implementation. You will not merely learn to use PyTorch—you will understand *why* it works the way it does, and you will develop the intuition to design your own architectures for scientific problems.

How to Use This Book:

- **Theory sections** provide mathematical foundations and intuition. Do not skip these—they are the bedrock upon which everything rests.
- **Implementation sections** show you the "PyTorch way"—idioms, best practices, and common pitfalls. Type out the code yourself; reading is not enough.
- **Exercises** range from elementary to research-level. Attempt all of them. Struggle is where learning happens.
- **Colored boxes** highlight key insights:
 - Professor's Notes for expert intuition
 - PyTorch Idioms for best practices
 - Warnings for common mistakes
 - Expert Tricks for advanced techniques
 - Scientific Applications showing real-world use

Treat this not just as a technical exercise, but as an art. The code we write, like the L^AT_EX that renders this page, can and should be clean, precise, and beautiful.

Let us begin.

Professor Vectorex Gradiens
ETH Zürich
December 2025

Contents

Preface	3
I Foundations of PyTorch	7
1 The torch.Tensor	7
1.1 Introduction: Why Tensors?	7
1.2 Theory: Tensor Fundamentals	7
1.3 Implementation: Creating and Manipulating Tensors	8
1.3.1 Creating Tensors	8
1.3.2 Tensor Properties	9
1.3.3 Indexing and Slicing	9
1.3.4 Shape Manipulation	9
1.3.5 Broadcasting	11
1.3.6 Mathematical Operations	11
1.3.7 Moving Tensors Between Devices	13
1.4 Exercises	13
2 Autograd: Automatic Differentiation	16
2.1 Introduction: Why Autograd Matters	16
2.2 Theory: Computational Graphs and Backpropagation	16
2.2.1 The Computational Graph	16
2.2.2 Backpropagation: The Chain Rule	16
2.2.3 Dynamic Computation Graphs	17
2.3 Implementation: Using Autograd in Practice	17
2.3.1 Basic Gradient Computation	17
2.3.2 Leaf vs Non-Leaf Tensors	17
2.3.3 Gradient Accumulation and Zeroing	19
2.3.4 Vector-Jacobian Products	19
2.3.5 Disabling Gradient Tracking	20
2.3.6 In-Place Operations	21
2.3.7 Retaining Graphs	21
2.3.8 Higher-Order Derivatives	21
2.4 Implementation: Common Patterns and Debugging	22
2.4.1 Checking Gradients	22
2.4.2 Gradient Flow Debugging	23
2.4.3 Custom Autograd Functions	24
2.5 Exercises	24
3 Building Models with nn.Module	27
3.1 Introduction: Why nn.Module?	27
3.2 Theory: The Module Hierarchy	27
3.3 Implementation: Creating Your First Module	27
3.3.1 A Simple Linear Model	27
3.3.2 Accessing Parameters	28
3.3.3 Using nn.Sequential	28
3.3.4 Custom Modules with Multiple Layers	30
3.3.5 Parameters vs Buffers	30
3.3.6 Training vs Evaluation Mode	31
3.3.7 Moving Models to GPU	32
3.3.8 Saving and Loading Models	32
3.3.9 Parameter Initialization	33

3.3.10	Module Hooks for Debugging	35
3.4	Implementation: Practical Patterns	35
3.4.1	Freezing Layers	35
3.4.2	Multiple Inputs/Outputs	36
3.4.3	Custom Layer with Parameters	36
3.5	Exercises	37
4	The Training Loop & Optimizers	40
4.1	Introduction: The Training Process	40
4.2	Theory: Gradient Descent and Optimization	40
4.2.1	The Optimization Problem	40
4.2.2	Gradient Descent	40
4.2.3	Learning Rate	40
4.3	Implementation: The Standard Training Loop	41
4.3.1	Complete Training Example	41
4.3.2	Loss Functions	43
4.3.3	Optimizers	44
4.3.4	Learning Rate Schedules	45
4.3.5	Gradient Clipping	45
4.3.6	Validation Loop	47
4.3.7	Early Stopping	48
4.4	Implementation: Debugging Training	49
4.4.1	When Training Goes Wrong	49
4.4.2	Monitoring Training	50
4.5	Exercises	51
5	Data Loading & Preprocessing	54
5.1	Introduction: Why Proper Data Loading Matters	54
5.2	Theory: Dataset and DataLoader Architecture	54
5.3	Implementation: Basic Usage	54
5.3.1	Using Built-in Datasets	54
5.3.2	Creating Custom Datasets	56
5.3.3	Dataset from Files	56
5.3.4	Dataset for Time Series	57
5.3.5	Dataset for Point Clouds	59
5.3.6	Data Transforms	59
5.3.7	Normalization Strategies	61
5.3.8	Handling Variable-Length Sequences	62
5.3.9	Train/Validation/Test Split	64
5.3.10	Efficient Data Loading	64
5.4	Exercises	66
II	Deep Learning Architectures	68
6	Fully Connected Networks (MLPs)	68
6.1	Introduction: The Foundation of Deep Learning	68
6.2	Theory: Understanding MLPs	68
6.2.1	Architecture	68
6.2.2	Depth vs Width	69
6.2.3	Activation Functions	70
6.2.4	Why Networks Fail to Train	72
6.3	Implementation: Building MLPs in PyTorch	73
6.3.1	Simple MLP with Sequential	73
6.3.2	Custom MLP Module	73

6.3.3	Weight Initialization	75
6.3.4	Batch Normalization	77
6.3.5	Layer Normalization	78
6.3.6	Dropout	79
6.3.7	Complete MLP with All Techniques	79
6.4	Implementation: Debugging MLPs	81
6.4.1	Checking for Dead ReLUs	81
6.4.2	Visualizing Activation Distributions	81
6.4.3	Gradient Flow Visualization	82
6.5	Exercises	84
6.6	Key Takeaways	87
7	Convolutional Neural Networks (CNNs)	88
7.1	Introduction: Why Convolutions?	88
7.2	Theory: How Convolutions Work	88
7.2.1	The Convolution Operation	88
7.2.2	Parameter Sharing	89
7.2.3	Receptive Fields	90
7.2.4	Output Size Calculation	90
7.2.5	Pooling	91
7.3	Implementation: Building CNNs in PyTorch	92
7.3.1	Conv2d: The Core Layer	92
7.3.2	Shape Calculation Practice	92
7.3.3	Building a Simple CNN	94
7.3.4	Modern CNN with Batch Normalization	94
7.3.5	1D Convolutions for Time Series	97
7.3.6	3D Convolutions for Volumetric Data	97
7.3.7	Transposed Convolutions (Upsampling)	99
7.3.8	Dilated Convolutions	100
7.3.9	Common CNN Architectures	101
7.4	Implementation: Debugging CNNs	101
7.4.1	Visualizing Feature Maps	101
7.4.2	Common CNN Debugging Issues	102
7.5	Exercises	104
7.6	Key Takeaways	107
8	Residual Networks & Skip Connections	108
8.1	Introduction: The Deep Network Problem	108
8.2	Theory: Why Skip Connections Work	108
8.2.1	The Problem: Vanishing Gradients	108
8.2.2	The Solution: Residual Connections	108
8.2.3	Gradient Flow Through Skip Connections	109
8.2.4	Residual Block Variants	110
8.2.5	When Dimensions Don't Match	111
8.3	Implementation: Building Residual Blocks	111
8.3.1	Basic Residual Block	111
8.3.2	Pre-Activation Residual Block	113
8.3.3	Bottleneck Block	113
8.3.4	Building a Complete ResNet	116
8.3.5	Dense Connections (DenseNet)	118
8.3.6	Skip Connections in Practice	119
8.4	Implementation: Debugging ResNets	120
8.4.1	Verifying Skip Connections	120
8.4.2	Comparing With and Without Skip Connections	120
8.5	Exercises	121

8.6	Key Takeaways	124
9	Batch Normalization & Layer Normalization	126
9.1	Introduction: Why Normalization Matters	126
9.2	Theory: Batch Normalization	126
9.2.1	The Problem: Internal Covariate Shift	126
9.2.2	How Batch Normalization Works	126
9.2.3	Why Batch Normalization Helps	128
9.2.4	Limitations of Batch Normalization	128
9.3	Theory: Layer Normalization	128
9.4	Implementation: Using Normalization in PyTorch	130
9.4.1	Batch Normalization Usage	130
9.4.2	Layer Normalization Usage	131
9.4.3	Other Normalization Variants	132
9.4.4	Implementing Batch Norm from Scratch	132
9.4.5	Common Mistakes and Debugging	134
9.5	Exercises	136
9.6	Key Takeaways	139
10	Recurrent Neural Networks (RNNs)	141
10.1	Introduction: Processing Sequential Data	141
10.2	Theory: How RNNs Work	141
10.2.1	The Basic RNN	141
10.2.2	Parameter Sharing	143
10.2.3	Sequence Modeling Patterns	143
10.2.4	The Vanishing Gradient Problem	143
10.3	Theory: Long Short-Term Memory (LSTM)	144
10.3.1	LSTM Architecture	144
10.4	Theory: Gated Recurrent Unit (GRU)	146
10.4.1	GRU Architecture	146
10.5	Implementation: RNNs in PyTorch	146
10.5.1	Vanilla RNN	146
10.5.2	LSTM	148
10.5.3	GRU	148
10.5.4	Building a Sequence Classifier	148
10.5.5	Sequence-to-Sequence (Many-to-Many)	150
10.5.6	Bidirectional RNNs	150
10.5.7	Handling Variable-Length Sequences	152
10.5.8	Gradient Clipping (Essential for RNNs)	153
10.5.9	Training Tips for RNNs	153
10.5.10	Common RNN Architectures	155
10.5.11	Time Series Forecasting Example	157
10.6	Exercises	158
10.7	Key Takeaways	162
11	Attention & Transformers	164
11.1	Introduction: Beyond RNNs	164
11.2	Theory: The Attention Mechanism	164
11.2.1	Intuition: What is Attention?	164
11.2.2	Mathematical Formulation	164
11.2.3	Self-Attention vs Cross-Attention	166
11.2.4	Multi-Head Attention	166
11.2.5	Positional Encoding	166
11.3	Theory: The Transformer Architecture	168
11.3.1	Complete Architecture	168

11.3.2	Why Transformers Work So Well	168
11.4	Implementation: Building Attention from Scratch	170
11.4.1	Scaled Dot-Product Attention	170
11.4.2	Multi-Head Attention	171
11.4.3	Positional Encoding	174
11.4.4	Feed-Forward Network	175
11.4.5	Complete Transformer Encoder Block	175
11.4.6	Using PyTorch's Built-in Attention	177
11.4.7	Masking for Causal Attention	177
11.4.8	Complete Transformer for Classification	178
11.4.9	Transformer for Sequence-to-Sequence	180
11.5	Exercises	183
11.6	Key Takeaways	188
III	Practical Skills	190
12	Debugging & Best Practices	190
12.1	Introduction: Why Debugging Matters	190
12.2	Common Bugs and How to Fix Them	190
12.2.1	Bug 1: Loss is NaN	190
12.2.2	Bug 2: Loss Not Decreasing	192
12.2.3	Bug 3: Training Loss Decreases, Validation Doesn't	194
12.2.4	Bug 4: Model Works on Small Data, Fails on Full Dataset	195
12.3	Debugging Strategies	196
12.3.1	Start Simple, Add Complexity	196
12.3.2	Check Shapes at Every Step	196
12.3.3	Visualize Gradients	197
12.3.4	Use Hooks for Debugging	197
12.4	Best Practices for Development	198
12.4.1	Code Organization	198
12.4.2	Configuration Management	199
12.4.3	Experiment Tracking	199
12.4.4	Reproducibility	201
12.4.5	Checkpointing	201
12.5	Debugging Checklist	202
12.6	Key Takeaways	203
13	Training Best Practices	204
13.1	Introduction: From Working to Working Well	204
13.2	Hyperparameter Tuning	204
13.2.1	Which Hyperparameters Matter Most?	204
13.2.2	Learning Rate Strategies	204
13.2.3	Batch Size Selection	207
13.2.4	Model Architecture Selection	207
13.3	Regularization Strategies	208
13.3.1	Dropout	208
13.3.2	Weight Decay (L2 Regularization)	208
13.3.3	Data Augmentation	208
13.4	Training Monitoring	209
13.4.1	What to Track	209
13.4.2	Recognizing Training Issues from Curves	210
13.5	Model Selection and Evaluation	211
13.5.1	Proper Train/Val/Test Split	211
13.5.2	Cross-Validation for Small Datasets	211

13.6	Practical Tips	213
13.6.1	Start with Good Defaults	213
13.6.2	Hyperparameter Search Strategies	213
13.7	Key Takeaways	214
14	Performance & Optimization	216
14.1	Introduction: Making Training Faster	216
14.2	Mixed Precision Training	216
14.2.1	Automatic Mixed Precision (AMP)	216
14.2.2	Gradient Checkpointing	217
14.3	GPU Optimization	218
14.3.1	Data Transfer Optimization	218
14.3.2	Batch Size and GPU Utilization	218
14.3.3	Multi-GPU Training	218
14.4	DataLoader Optimization	220
14.4.1	Optimal Number of Workers	220
14.4.2	Prefetching	220
14.5	Memory Management	220
14.5.1	Monitoring Memory Usage	220
14.5.2	Reducing Memory Usage	221
14.6	Profiling and Bottleneck Analysis	223
14.6.1	PyTorch Profiler	223
14.6.2	Simple Timing	223
14.6.3	Identify Bottlenecks	223
14.7	Production Deployment	226
14.7.1	Model Export	226
14.7.2	Inference Optimization	226
14.8	Key Takeaways	227
IV	Integrative Challenges	229
15	Grand Challenges	229
15.1	Introduction: Putting It All Together	229
15.2	Challenge 1: Time Series Forecasting System	229
15.3	Challenge 2: Hybrid Vision-Language Model	231
15.4	Challenge 3: Denoising Scientific Data	234
15.5	Challenge 4: Custom Transformer for Sequence Tasks	238
15.6	Challenge 5: End-to-End ML Pipeline	239
15.7	Evaluation and Next Steps	240
15.8	Final Thoughts	241

Part I

Foundations of PyTorch

1 The torch.Tensor

1.1 Introduction: Why Tensors?

Everything in PyTorch—from your input data (images, time series, simulation meshes) to the parameters of your neural network (weights and biases)—is represented as a **tensor**.

If you’ve used NumPy, you already know the core idea: a tensor is essentially a multi-dimensional array (like `ndarray`). So why use PyTorch instead of NumPy?

Two critical reasons:

1. **GPU Acceleration:** PyTorch tensors can be moved to GPUs with a single command, enabling massive parallelization. For deep learning and large-scale scientific computing, this isn’t optional—it’s essential.
2. **Automatic Differentiation:** PyTorch can automatically compute gradients of any output with respect to any input. This **autograd** system is the engine that makes deep learning possible. We’ll explore this in detail in Section 2.

1.2 Theory: Tensor Fundamentals

Tensor

A **tensor** is a multi-dimensional array of numerical values. It generalizes:

- **Scalar** (0D tensor): a single number
- **Vector** (1D tensor): an array of numbers
- **Matrix** (2D tensor): a 2D grid of numbers
- **Higher-dimensional tensors:** 3D, 4D, ... arrays

Key Properties:

- **Shape:** The dimensions of the tensor, e.g., (3, 4, 5) means $3 \times 4 \times 5$
- **Dtype:** The data type (e.g., `float32`, `int64`, `bool`)
- **Device:** Where the tensor lives (`cpu` or `cuda:0`, `cuda:1`, etc.)
- **Requires_grad:** Whether to track operations for automatic differentiation

Common Tensor Shapes in Deep Learning:

Shape	Typical Use
(N,)	Vector of N values
(N, D)	Batch of N samples, D features each
(N, C, H, W)	Batch of N images: C channels, H×W pixels
(N, L, D)	Batch of N sequences, length L, D features
(N, C, D, H, W)	Batch of N 3D volumes

Table 1: Common tensor shapes. N=batch size, C=channels, H/W=height/width, L=sequence length, D=feature dimension

1.3 Implementation: Creating and Manipulating Tensors

1.3.1 Creating Tensors

```
1 import torch
2
3 # From Python lists
4 x = torch.tensor([1, 2, 3])
5 print(x) # tensor([1, 2, 3])
6
7 # Specify dtype
8 x = torch.tensor([1.0, 2.0, 3.0], dtype=torch.float32)
9 print(x.dtype) # torch.float32
10
11 # Common initialization functions
12 zeros = torch.zeros(2, 3) # 2x3 matrix of zeros
13 ones = torch.ones(2, 3) # 2x3 matrix of ones
14 random = torch.randn(2, 3) # Random normal
15 # distribution
16 uniform = torch.rand(2, 3) # Uniform distribution
17 # [0, 1)
18
19 arange = torch.arange(0, 10, 2) # [0, 2, 4, 6, 8]
20 linspace = torch.linspace(0, 1, 5) # [0.0, 0.25, 0.5, 0.75,
21 # 1.0]
22
23 # Like existing tensor (same shape/dtype)
24 x = torch.ones(3, 4)
25 y = torch.zeros_like(x) # Same shape as x
26 z = torch.randn_like(x) # Same shape, random values
```

[Default Dtype] PyTorch’s default floating point dtype is `torch.float32`. For most deep learning, this is fine. You can change it globally with:

```
1 torch.set_default_dtype(torch.float64) # More precision
```

However, `float32` is usually the sweet spot: enough precision, less memory, faster on GPUs.

1.3.2 Tensor Properties

```

1 x = torch.randn(2, 3, 4)
2
3 print(x.shape)           # torch.Size([2, 3, 4])
4 print(x.size())          # Same as .shape
5 print(x.ndim)            # 3 (number of dimensions)
6 print(x.numel())         # 24 (total number of elements)
7 print(x.dtype)           # torch.float32
8 print(x.device)          # cpu
9
10 # Check if tensor requires gradients
11 print(x.requires_grad)  # False by default

```

1.3.3 Indexing and Slicing

PyTorch indexing works like NumPy:

```

1 x = torch.arange(12).reshape(3, 4)
2 # tensor([[ 0,  1,  2,  3],
3 #         [ 4,  5,  6,  7],
4 #         [ 8,  9, 10, 11]])
5
6 # Basic indexing
7 print(x[0])              # First row: tensor([0, 1, 2, 3])
8 print(x[:, 0])           # First column: tensor([0, 4, 8])
9 print(x[1, 2])           # Element at (1, 2): tensor(6)
10
11 # Slicing
12 print(x[:2, :2])        # Top-left 2x2:
13 # tensor([[0, 1],
14 #         [4, 5]])
15
16 # Advanced indexing
17 indices = torch.tensor([0, 2])
18 print(x[indices])       # Rows 0 and 2
19
20 # Boolean masking
21 mask = x > 5
22 print(x[mask])          # tensor([ 6,  7,  8,  9, 10, 11])

```

[Views vs Copies] Slicing creates a **view**, not a copy! Modifying a view changes the original:

```

1 x = torch.ones(3, 3)
2 y = x[0]                # View of first row
3 y[0] = 999              # Modifies x!
4 print(x)                # First element is now 999
5
6 # To create a copy:
7 y = x[0].clone()

```

1.3.4 Shape Manipulation

```

1 x = torch.arange(12)
2
3 # Reshape (must have same total elements)
4 y = x.reshape(3, 4)
5 z = x.view(2, 6)        # Similar to reshape

```

```
6 w = x.reshape(3, -1) # -1 means "infer this dimension" ->
   (3, 4)
7
8 # Transpose
9 x = torch.randn(2, 3)
10 y = x.T # Shape: (3, 2)
11 z = x.transpose(0, 1) # Same as .T
12
13 # Permute (generalized transpose)
14 x = torch.randn(2, 3, 4)
15 y = x.permute(2, 0, 1) # Shape: (4, 2, 3)
16
17 # Flatten
18 x = torch.randn(2, 3, 4)
19 y = x.flatten() # Shape: (24,)
20 z = x.flatten(1) # Flatten from dim 1: (2, 12)
21
22 # Squeeze and unsqueeze (remove/add dimensions of size 1)
23 x = torch.randn(2, 1, 3, 1)
24 y = x.squeeze() # Shape: (2, 3)
25 z = x.squeeze(1) # Remove dim 1: (2, 3, 1)
26
27 x = torch.randn(2, 3)
28 y = x.unsqueeze(0) # Shape: (1, 2, 3)
29 z = x.unsqueeze(-1) # Shape: (2, 3, 1)
```

[Understanding Dimensions]

- Dimension indices start at 0
- Negative indices count from the end: -1 is the last dimension
- `unsqueeze(1)` adds a dimension at position 1
- `squeeze()` removes ALL dimensions of size 1

1.3.5 Broadcasting

Broadcasting is one of the most powerful—and confusing—features. It allows operations between tensors of different shapes.

Broadcasting Rules:

1. If tensors have different number of dimensions, prepend 1s to the smaller one
2. Dimensions are compatible if they're equal or one of them is 1
3. The result shape is the maximum along each dimension

```

1 # Example 1: Vector + Scalar
2 x = torch.tensor([1, 2, 3]) # Shape: (3,)
3 y = 10                      # Shape: () - scalar
4 z = x + y                   # Shape: (3,) - broadcasts
                             # scalar to (3,)
5 # Result: tensor([11, 12, 13])
6
7 # Example 2: Matrix + Vector (row-wise)
8 x = torch.ones(3, 4)        # Shape: (3, 4)
9 y = torch.tensor([1, 2, 3, 4]) # Shape: (4,)
10 z = x + y                   # y broadcasts to (1, 4) then
                             # (3, 4)
11 # Each row of x has y added to it
12
13 # Example 3: Matrix + Column Vector
14 x = torch.ones(3, 4)        # Shape: (3, 4)
15 y = torch.tensor([[1], [2], [3]]) # Shape: (3, 1)
16 z = x + y                   # y broadcasts to (3, 4)
17 # Each column of x has corresponding y value added
18
19 # Example 4: Batch operations
20 x = torch.randn(10, 3, 4)    # 10 matrices of size 3x4
21 y = torch.randn(3, 4)        # Single matrix
22 z = x + y                   # y broadcasts to (10, 3, 4)
23 # Same matrix y added to all 10 matrices in x

```

[Common Broadcasting Mistake]

```

1 # Trying to add vectors of different sizes
2 x = torch.ones(3)           # Shape: (3,)
3 y = torch.ones(4)           # Shape: (4,)
4 z = x + y                   # ERROR! Shapes don't match
5
6 # Need to explicitly reshape:
7 x = x.unsqueeze(1)          # Shape: (3, 1)
8 y = y.unsqueeze(0)          # Shape: (1, 4)
9 z = x + y                   # Shape: (3, 4) - outer sum

```

1.3.6 Mathematical Operations

```

1 x = torch.tensor([1.0, 2.0, 3.0])
2 y = torch.tensor([4.0, 5.0, 6.0])
3
4 # Element-wise operations
5 z = x + y                   # Addition
6 z = x - y                   # Subtraction
7 z = x * y                   # Multiplication
8 z = x / y                   # Division
9 z = x ** 2                  # Power
10 z = torch.sqrt(x)          # Square root

```

```

11 z = torch.exp(x)      # Exponential
12 z = torch.log(x)      # Natural log
13
14 # In-place operations (modify the tensor)
15 x.add_(5)              # Add 5 to x in-place
16 x.mul_(2)              # Multiply x by 2 in-place
17
18 # Matrix operations
19 A = torch.randn(3, 4)
20 B = torch.randn(4, 5)
21 C = torch.mm(A, B)      # Matrix multiply: (3, 5)
22 C = torch.matmul(A, B)  # Same, but works for batches
23 C = A @ B              # Python 3.5+ operator
24
25 # Batch matrix multiply
26 A = torch.randn(10, 3, 4) # 10 matrices of size 3x4
27 B = torch.randn(10, 4, 5) # 10 matrices of size 4x5
28 C = torch.bmm(A, B)      # Shape: (10, 3, 5)
29
30 # Reduction operations
31 x = torch.randn(3, 4)
32 mean = x.mean()          # Scalar: mean of all elements
33 std = x.std()            # Standard deviation
34 sum_all = x.sum()        # Sum of all elements
35 sum_dim0 = x.sum(dim=0)   # Sum along dimension 0: (4,)
36 sum_dim1 = x.sum(dim=1)   # Sum along dimension 1: (3,)
37 max_val, max_idx = x.max(dim=1) # Max value and index per
    row

```

[Dimension Conventions] When reducing along a dimension:

- `dim=0`: Operate across rows (result has fewer rows)
- `dim=1`: Operate across columns (result has fewer columns)
- `dim=-1`: Last dimension (most common for sequences)
- `keepdim=True`: Keep the dimension (size 1) in result

1.3.7 Moving Tensors Between Devices

```

1 # Check if CUDA is available
2 if torch.cuda.is_available():
3     device = torch.device("cuda")
4     print(f"Using GPU: {torch.cuda.get_device_name(0)}")
5 else:
6     device = torch.device("cpu")
7     print("Using CPU")
8
9 # Create tensor on CPU
10 x = torch.randn(1000, 1000)
11
12 # Move to GPU
13 x_gpu = x.to(device)
14 # Or:
15 x_gpu = x.cuda() # Shorthand for .to('cuda')
16
17 # Move back to CPU
18 x_cpu = x_gpu.to('cpu')
19 # Or:
20 x_cpu = x_gpu.cpu()
21
22 # Operations must be on same device
23 y_gpu = torch.randn(1000, 1000, device=device) # Create
    directly on GPU
24 z = x_gpu + y_gpu # OK: both on GPU
25
26 # This would ERROR:
27 # z = x_cpu + y_gpu # ERROR! Different devices

```

[Device Mismatch] One of the most common errors:

```
1 RuntimeError: Expected all tensors to be on the same device
```

Always ensure tensors are on the same device before operations! Use `.to(device)` consistently.

1.4 Exercises

1.1: Tensor Creation - ★★

Goal: Get comfortable creating tensors in different ways.

1. Create a tensor of shape (5, 3) filled with random values from a standard normal distribution
2. Create a tensor of shape (4, 4) filled with values from 0 to 15 (use `arange` and `reshape`)
3. Create a tensor of shape (10,) with evenly spaced values from 0 to 1
4. Create a (3, 3) identity matrix (hint: `torch.eye`)

Starter code:

```

1 import torch
2
3 # Your code here
4 x1 = ...
5 x2 = ...
6 x3 = ...
7 x4 = ...
8

```



```
9 print(x1.shape, x2.shape, x3.shape, x4.shape)
```

1.2: Shape Manipulation - ★★

Goal: Master reshaping and dimension manipulation.

Given `x = torch.arange(24)`:

1. Reshape `x` to (2, 3, 4)
2. Transpose the last two dimensions to get (2, 4, 3)
3. Flatten the result to 1D
4. Create a view with an extra dimension at position 0 to get (1, 24)

Hint: Use `reshape`, `transpose`, `flatten`, `unsqueeze`.

1.3: Broadcasting Mastery - ★★★

Goal: Understand broadcasting and avoid common mistakes.

1. Create matrix `A` of shape (5, 1) with values [1, 2, 3, 4, 5]
2. Create vector `b` of shape (3,) with values [10, 20, 30]
3. Compute `C = A + b` using broadcasting. What is the shape of `C`?
4. Compute the outer product of two vectors `u` (size 4) and `v` (size 3) without loops (result should be 4×3)
5. Given batch of images (16, 3, 32, 32) and per-channel mean (3,), subtract the mean from each image

Debugging tip: If broadcasting fails, print the shapes! Use `print(A.shape, b.shape)`.

1.4: Advanced Indexing - ★★★

Goal: Learn advanced indexing techniques.

Given `x = torch.randn(100, 5)`:

1. Select all rows where the first column is positive
2. Get the maximum value in each row and its index
3. Create a new tensor with only the even-indexed rows (0, 2, 4, ...)
4. Replace all negative values with 0 (use boolean indexing)

Hint: Use boolean masks: `mask = x[:, 0] > 0`, then `x[mask]`.

1.5: Implementing Distance Matrix - ★★★★★

Goal: Use broadcasting for efficient computation.

Given a set of points `X` of shape (N, D) where N=100 points in D=2 dimensions:

1. Compute the pairwise Euclidean distance matrix (N×N) **without loops**
2. The (i, j) entry should be $\|x_i - x_j\|_2$
3. Verify the diagonal is all zeros
4. Find the nearest neighbor for each point (excluding itself)

Hint: Expand dimensions: `X_i = X.unsqueeze(1)` (shape N, 1, D) and `X_j = X.unsqueeze(0)` (shape 1, N, D). Then compute `(X_i - X_j)**2`.

1.6: GPU Operations - ★★★

Goal: Practice moving tensors between devices.

1. Check if CUDA is available on your system
2. Create a large tensor (1000×1000) on CPU
3. Move it to GPU (if available) and time a matrix multiplication
4. Compare the time with the same operation on CPU
5. Handle the case where GPU is not available gracefully

Starter code:

```
1 import torch
2 import time
3
4 device = torch.device("cuda" if torch.cuda.is_available
    () else "cpu")
5 print(f"Using device: {device}")
6
7 # Your code here
```

2 Autograd: Automatic Differentiation

2.1 Introduction: Why Autograd Matters

Deep learning is fundamentally an optimization problem. We have a model with parameters (weights), a loss function that measures how wrong our predictions are, and we want to adjust the parameters to minimize the loss.

The key insight: we use **gradient descent**. We compute the gradient of the loss with respect to each parameter, then update the parameters in the direction that decreases the loss:

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla_{\theta} \mathcal{L}$$

For a neural network with millions of parameters, computing these gradients by hand is impossible. PyTorch's **autograd** system does this automatically through the chain rule.

2.2 Theory: Computational Graphs and Backpropagation

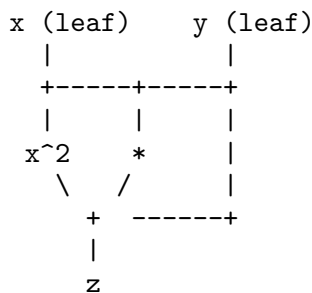
2.2.1 The Computational Graph

Every operation you perform on a tensor creates a **computational graph**. This graph tracks how each tensor was computed from other tensors.

Example: Simple computation

```
1 import torch
2
3 x = torch.tensor(2.0, requires_grad=True)
4 y = torch.tensor(3.0, requires_grad=True)
5 z = x * y + x**2  # z = xy + x^2
```

This creates a graph:



Key Concepts:

- **Leaf tensors:** Input tensors created by the user with `requires_grad=True`
- **Non-leaf tensors:** Result of operations on other tensors
- **grad_fn:** Each non-leaf tensor remembers the operation that created it

2.2.2 Backpropagation: The Chain Rule

When you call `z.backward()`, PyTorch:

1. Traverses the graph backwards from z
2. Applies the chain rule at each operation

3. Accumulates gradients in the `.grad` attribute of leaf tensors

For our example:

$$\frac{\partial z}{\partial x} = \frac{\partial}{\partial x}(xy + x^2) = y + 2x$$

$$\frac{\partial z}{\partial y} = \frac{\partial}{\partial y}(xy + x^2) = x$$

With $x = 2, y = 3$: $\frac{\partial z}{\partial x} = 7, \frac{\partial z}{\partial y} = 2$

[Why "Automatic"?] You never write the derivative formulas! PyTorch knows the derivative of every operation (add, multiply, sin, exp, matrix multiply, etc.) and chains them automatically. This is the magic of autograd.

2.2.3 Dynamic Computation Graphs

PyTorch uses **dynamic computation graphs** (define-by-run):

- The graph is built *during* the forward pass
- You can use Python control flow (if statements, loops)
- The graph can be different each time you run the code
- After `.backward()`, the graph is *destroyed* (by default)

This is different from TensorFlow 1.x (static graphs) and makes PyTorch more flexible for research.

2.3 Implementation: Using Autograd in Practice

2.3.1 Basic Gradient Computation

```

1 import torch
2
3 # Create tensors with gradient tracking
4 x = torch.tensor(2.0, requires_grad=True)
5 y = torch.tensor(3.0, requires_grad=True)
6
7 # Forward pass: compute z
8 z = x**2 + 2*x*y + y**2
9
10 print(z) # tensor(25., grad_fn=<AddBackward0>)
11 print(z.grad_fn) # Shows the operation that created z
12
13 # Backward pass: compute gradients
14 z.backward()
15
16 # Gradients are now in .grad
17 print(x.grad) # dz/dx = 2x + 2y = 10
18 print(y.grad) # dz/dy = 2x + 2y = 10

```

2.3.2 Leaf vs Non-Leaf Tensors

```

1 x = torch.tensor(2.0, requires_grad=True)
2 y = x * 2
3
4 print(x.is_leaf) # True - created by user
5 print(y.is_leaf) # False - result of operation
6
7 # Only leaf tensors store gradients by default

```

```
8 z = y.sum()
9 z.backward()
10
11 print(x.grad)    # Available! x is a leaf
12 print(y.grad)    # None! Non-leaf gradients not retained
```

To retain gradients for non-leaf tensors:

```
1 x = torch.tensor([1.0, 2.0, 3.0], requires_grad=True)
2 y = x * 2
3 y.retain_grad()    # Tell PyTorch to keep this gradient
4
5 z = y.sum()
6 z.backward()
7
8 print(x.grad)    # tensor([2., 2., 2.])
9 print(y.grad)    # tensor([1., 1., 1.]) - now available!
```

2.3.3 Gradient Accumulation and Zeroing

[Gradients Accumulate!] Calling `.backward()` multiple times **adds** to existing gradients. You must zero them manually between iterations.

```

1 x = torch.tensor(2.0, requires_grad=True)
2
3 # First backward pass
4 y = x**2
5 y.backward()
6 print(x.grad)  # tensor(4.)
7
8 # Second backward pass without zeroing
9 y = x**2
10 y.backward()
11 print(x.grad)  # tensor(8.) - DOUBLED!
12
13 # Correct approach: zero gradients first
14 x.grad.zero_()
15 y = x**2
16 y.backward()
17 print(x.grad)  # tensor(4.) - correct again

```

Why does this matter? In training loops, we compute gradients for each batch and update weights. Between batches, we must zero gradients:

```

1 # Training loop pattern (simplified)
2 for batch in dataloader:
3     # Zero gradients from previous batch
4     optimizer.zero_grad()  # Or: model.zero_grad()
5
6     # Forward pass
7     output = model(batch)
8     loss = criterion(output, target)
9
10    # Backward pass
11    loss.backward()
12
13    # Update weights
14    optimizer.step()

```

2.3.4 Vector-Jacobian Products

`.backward()` can only be called on scalar outputs (or you must provide a gradient argument):

```

1 # Scalar output - works
2 x = torch.tensor([1.0, 2.0, 3.0], requires_grad=True)
3 y = x.sum()  # Scalar
4 y.backward()
5 print(x.grad)  # Works!
6
7 # Vector output - need gradient argument
8 x = torch.tensor([1.0, 2.0, 3.0], requires_grad=True)
9 y = x * 2  # Vector: [2, 4, 6]
10 # y.backward()  # ERROR! Output is not scalar
11
12 # Provide gradient (weights for each output)
13 gradient = torch.tensor([1.0, 1.0, 1.0])
14 y.backward(gradient)
15 print(x.grad)  # tensor([2., 2., 2.])

```

What's happening? When output is a vector \mathbf{y} , we compute:

$$\mathbf{x}.\text{grad} = \frac{\partial \mathbf{y}^T}{\partial \mathbf{x}} \cdot \text{gradient}$$

This is a **vector-Jacobian product**. For scalar loss, gradient is implicitly 1.

2.3.5 Disabling Gradient Tracking

Sometimes you don't need gradients (inference, debugging, certain computations):

```
1 x = torch.tensor([1.0, 2.0], requires_grad=True)
2
3 # Method 1: torch.no_grad() context
4 with torch.no_grad():
5     y = x * 2
6     z = y.sum()
7 print(z.requires_grad) # False
8 # z.backward() # ERROR - can't backprop
9
10 # Method 2: detach() - creates a view without gradients
11 y = x * 2
12 y_detached = y.detach()
13 print(y.requires_grad) # True
14 print(y_detached.requires_grad) # False
15
16 # Method 3: requires_grad_() - in-place toggle
17 x = torch.tensor([1.0, 2.0], requires_grad=True)
18 x.requires_grad_(False)
19 print(x.requires_grad) # False
```

[When to Disable Gradients]

- **Inference:** Use `torch.no_grad()` or `torch.inference_mode()` for evaluation
- **Freezing layers:** Set `param.requires_grad = False` for parameters you don't want to train
- **Saving memory:** Gradient tracking uses extra memory
- **Performance:** Operations are faster without gradient tracking

2.3.6 In-Place Operations

[In-Place Operations Break Gradients] In-place operations (those ending with `_`) can cause errors in backpropagation.

```

1 # BAD: In-place operation on tensor used in computation
2 x = torch.tensor([1.0, 2.0], requires_grad=True)
3 y = x * 2
4 x.add_(1) # Modifies x in-place
5 # y.backward(torch.ones_like(y)) # ERROR! x was modified
6
7 # GOOD: Don't modify tensors that are part of the graph
8 x = torch.tensor([1.0, 2.0], requires_grad=True)
9 y = x * 2
10 z = y.sum()
11 z.backward() # Works!

```

However, in-place operations are fine for parameters during optimization (outside the graph):

```

1 # This is OK - happens after .backward()
2 with torch.no_grad():
3     x.add_(-learning_rate * x.grad)

```

2.3.7 Retaining Graphs

By default, the computation graph is freed after `.backward()`:

```

1 x = torch.tensor(2.0, requires_grad=True)
2 y = x**2
3 y.backward()
4 # y.backward() # ERROR! Graph was freed
5
6 # To keep the graph:
7 x = torch.tensor(2.0, requires_grad=True)
8 y = x**2
9 y.backward(retain_graph=True)
10 y.backward() # OK now! But remember to zero gradients if
    needed

```

When is this useful? Multiple backward passes (rare), or when you need to call `.backward()` multiple times in a single forward pass.

2.3.8 Higher-Order Derivatives

PyTorch can compute gradients of gradients:

```

1 # First derivative
2 x = torch.tensor(2.0, requires_grad=True)
3 y = x**3 # y = x^3
4
5 dy_dx = torch.autograd.grad(y, x, create_graph=True)[0]
6 print(dy_dx) # 3*x^2 = 12
7
8 # Second derivative
9 d2y_dx2 = torch.autograd.grad(dy_dx, x)[0]
10 print(d2y_dx2) # 6*x = 12

```

Note: We use `torch.autograd.grad()` instead of `.backward()` here because it gives us the gradient as a tensor (useful for higher-order derivatives).

[torch.autograd.grad vs .backward()]

- `.backward()`: Accumulates gradients in `.grad` attribute. Standard for training.
- `torch.autograd.grad()`: Returns gradients as tensors. Useful for:
 - Higher-order derivatives
 - Multiple gradient computations
 - When you need gradient as a value, not accumulated

2.4 Implementation: Common Patterns and Debugging

2.4.1 Checking Gradients

When implementing custom operations or debugging, verify gradients numerically:

```

1 import torch
2
3 def numerical_gradient(f, x, eps=1e-5):
4     """Compute gradient numerically using finite differences.
5     """
6     grad = torch.zeros_like(x)
7     for i in range(x.numel()):
8         x_plus = x.clone()
9         x_plus.flatten()[i] += eps
10        x_minus = x.clone()
11        x_minus.flatten()[i] -= eps
12        grad.flatten()[i] = (f(x_plus) - f(x_minus)) / (2 *
13            eps)
14    return grad
15
16 # Test it
17 x = torch.tensor([2.0, 3.0], requires_grad=True)
18 f = lambda t: (t**2).sum()
19
20 # Analytical gradient
21 y = f(x)
22 y.backward()
23 analytical_grad = x.grad
24
25 # Numerical gradient
26 numerical_grad = numerical_gradient(f, x)
27
28 print("Analytical:", analytical_grad)    # [4., 6.]
29 print("Numerical:", numerical_grad)      # [4., 6.] (
30     approximately)
31 print("Close?", torch.allclose(analytical_grad,
32     numerical_grad))

```

PyTorch also provides `torch.autograd.gradcheck`:

```

1 from torch.autograd import gradcheck
2
3 x = torch.randn(3, 4, requires_grad=True, dtype=torch.float64)
4
5 func = lambda t: (t**2).sum()
6
7 # Check if gradients are correct
8 test = gradcheck(func, x, eps=1e-6)
9 print(f"Gradient check passed: {test}")

```

2.4.2 Gradient Flow Debugging

When training doesn't work, check if gradients are flowing:

```
1 import torch.nn as nn
2
3 # Simple model
4 model = nn.Sequential(
5     nn.Linear(10, 50),
6     nn.ReLU(),
7     nn.Linear(50, 1)
8 )
9
10 # Forward pass
11 x = torch.randn(32, 10)
12 y = model(x)
13 loss = y.mean()
14
15 # Backward pass
16 loss.backward()
17
18 # Check gradients
19 for name, param in model.named_parameters():
20     if param.grad is not None:
21         print(f"{name}: grad mean = {param.grad.abs().mean()
22             :.6f}")
23     else:
24         print(f"{name}: NO GRADIENT!")
```

Common issues:

- **All zeros:** Dead ReLUs, wrong loss function
- **None:** Part of model disconnected from loss
- **Very large:** Exploding gradients (lower learning rate, add gradient clipping)
- **Very small:** Vanishing gradients (different architecture, normalization)

2.4.3 Custom Autograd Functions

For operations PyTorch doesn't support, or to override gradients:

```

1 class MySquare(torch.autograd.Function):
2     @staticmethod
3     def forward(ctx, input):
4         """
5         Forward pass: compute output and save what's needed
6         for backward.
7         ctx is a context object for saving information.
8         """
9         ctx.save_for_backward(input)
10        return input ** 2
11
12    @staticmethod
13    def backward(ctx, grad_output):
14        """
15        Backward pass: compute gradient of loss w.r.t. input.
16        grad_output is the gradient of loss w.r.t. output.
17        Return gradient w.r.t. each input (input only here).
18        """
19        input, = ctx.saved_tensors
20        grad_input = grad_output * 2 * input
21        return grad_input
22
23    # Use it
24    x = torch.tensor([1.0, 2.0, 3.0], requires_grad=True)
25    y = MySquare.apply(x) # Note: .apply() not ()
26    z = y.sum()
27    z.backward()
28    print(x.grad) # [2., 4., 6.] = 2*x

```

We'll use this for the exercises!

2.5 Exercises

2.1: Basic Gradient Computation - ★★

Goal: Get comfortable with `.backward()` and `.grad`.

1. Create tensors $a = 3$, $b = 4$ with `requires_grad=True`
2. Compute $c = a^2 + b^2$
3. Call `c.backward()` and verify:
 - `a.grad = 2a = 6`
 - `b.grad = 2b = 8`
4. Now compute $d = 2a + 3b$ and call `d.backward()`
5. What are `a.grad` and `b.grad` now? Why?

Hint: Remember that gradients accumulate!

2.2: Matrix Gradients - ★★

Goal: Understand gradients for matrix operations.

1. Create matrix W of shape (3, 4) with random values, `requires_grad=True`
2. Create input x of shape (4,) with random values (no gradients needed)
3. Compute $y = Wx$ (matrix-vector product)

4. Compute $L = \|y\|^2$ (sum of squares)
5. Call `L.backward()` and check the shape of `W.grad`
6. Verify it has the same shape as W

Starter code:

```
1 import torch
2 torch.manual_seed(42)
3
4 W = torch.randn(3, 4, requires_grad=True)
5 x = torch.randn(4)
6
7 # Your code here
```

2.3: Gradient Accumulation - ★★★

Goal: Master gradient zeroing in training loops.

Implement a mini training loop (without a real model):

1. Create a parameter w initialized to 5.0
2. For 10 iterations:
 - Compute loss $= (w - 2)^2$ (we want w to reach 2)
 - Compute gradient
 - Update: $w = w - 0.1 \cdot \text{grad}$
 - **Zero the gradient**
3. Print the final value of w (should be close to 2)

Hint: Use `with torch.no_grad():` for the update step.

2.4: Gradient Detaching - ★★★

Goal: Understand when and why to detach gradients.

1. Create $x = [1, 2, 3]$ with gradients
2. Compute $y = x^2$
3. Detach y to create y_{detach}
4. Compute $z_1 = y.\text{sum}()$ and $z_2 = y_{\text{detach}}.\text{sum}()$
5. Try to call `z1.backward()` - what happens?
6. Try to call `z2.backward()` - what happens?
7. Why would detaching be useful in training?

2.5: Higher-Order Derivatives - ★★★

Goal: Compute second derivatives (useful for some scientific ML methods).

1. Create $x = 2.0$ with gradients
2. Compute $y = \sin(x)$
3. Use `torch.autograd.grad()` to compute $\frac{dy}{dx}$ with `create_graph=True`
4. Compute the second derivative $\frac{d^2y}{dx^2}$
5. Verify analytically: $y = \sin(x) \Rightarrow y' = \cos(x) \Rightarrow y'' = -\sin(x)$
6. At $x = 2$: $y'' = -\sin(2) \approx -0.909$

Starter code:

```
1 import torch
```

```

2
3 x = torch.tensor(2.0, requires_grad=True)
4 y = torch.sin(x)
5
6 # First derivative
7 dy_dx = torch.autograd.grad(y, x, create_graph=True)[0]
8
9 # Your code for second derivative

```

2.6: Custom Autograd Function - ★★★★★

Goal: Implement a custom operation with manual gradient.

Implement a custom **sigmoid** function:

1. Forward: $\sigma(x) = \frac{1}{1+e^{-x}}$
2. Backward: $\frac{d\sigma}{dx} = \sigma(x)(1 - \sigma(x))$
3. Inherit from `torch.autograd.Function`
4. Implement `forward()` and `backward()` static methods
5. Test it on $x = [0, 1, 2]$
6. Verify gradient using `torch.autograd.gradcheck`

Hint: In `backward()`, you'll need the output of `sigmoid`. You can either:

- Save the output in `forward` using `ctx.save_for_backward`
- Or recompute it from the saved input

Starter code:

```

1 class MySigmoid(torch.autograd.Function):
2     @staticmethod
3     def forward(ctx, input):
4         # Your code here
5         pass
6
7     @staticmethod
8     def backward(ctx, grad_output):
9         # Your code here
10        pass
11
12 # Test
13 x = torch.tensor([0.0, 1.0, 2.0], requires_grad=True)
14 y = MySigmoid.apply(x)

```

3 Building Models with `nn.Module`

3.1 Introduction: Why `nn.Module`?

In Section 2, we manually created tensors with `requires_grad=True` and computed gradients. While this works for simple cases, real neural networks have hundreds of layers and millions of parameters. We need a way to organize this complexity.

`nn.Module` is PyTorch’s base class for all neural network components. It provides:

- **Automatic parameter management:** Tracks all learnable parameters
- **GPU/CPU transfer:** Move entire models between devices with one command
- **Training/eval modes:** Control behavior of layers like dropout and batch norm
- **State management:** Save and load model weights
- **Composability:** Build complex models from simple modules

Every neural network layer (`nn.Linear`, `nn.Conv2d`, etc.) and every model you build inherits from `nn.Module`.

3.2 Theory: The Module Hierarchy

A `nn.Module` is a container that can hold:

- **Parameters:** Learnable tensors (weights, biases)
- **Buffers:** Non-learnable tensors (running statistics, constants)
- **Sub-modules:** Other `nn.Module` objects

Example hierarchy:

```
MyModel (nn.Module)
  layer1 (nn.Linear)
    weight (Parameter)
    bias (Parameter)
  layer2 (nn.Linear)
    weight (Parameter)
    bias (Parameter)
  activation (nn.ReLU)
```

When you call `model.parameters()`, PyTorch automatically finds all parameters in the entire tree.

[Key Principle: The Forward Method] Every `nn.Module` must implement a `forward()` method that defines the computation. You never call `forward()` directly—instead, you call the module as a function: `output = model(input)`. This triggers PyTorch’s hooks and tracking mechanisms.

3.3 Implementation: Creating Your First Module

3.3.1 A Simple Linear Model

Let’s build a simple linear regression model from scratch:

```
1 import torch
2 import torch.nn as nn
3
4 class LinearRegression(nn.Module):
5     def __init__(self, input_dim, output_dim):
6         """
7         Initialize the module.
```

```

8         Always call super().__init__() first!
9         """
10        super().__init__() # Initialize parent class
11
12        # Define learnable parameters
13        self.linear = nn.Linear(input_dim, output_dim)
14
15        def forward(self, x):
16            """
17            Define the forward pass.
18            x: input tensor of shape (batch_size, input_dim)
19            returns: output tensor of shape (batch_size,
20            output_dim)
21            """
22            return self.linear(x)
23
24        # Create the model
25        model = LinearRegression(input_dim=10, output_dim=1)
26
27        # Use the model
28        x = torch.randn(32, 10) # Batch of 32 samples
29        y = model(x) # Call the model (invokes forward())
30        print(y.shape) # torch.Size([32, 1])

```

[Always Call `super().__init__()`] Forgetting `super().__init__()` is a common mistake. Without it, PyTorch can't track parameters properly. Always put it as the first line of `__init__()`.

3.3.2 Accessing Parameters

```

1 model = LinearRegression(10, 1)
2
3 # Get all parameters
4 for name, param in model.named_parameters():
5     print(f"{name}: {param.shape}")
6
7 # Output:
8 # linear.weight: torch.Size([1, 10])
9 # linear.bias: torch.Size([1])
10
11 # Count parameters
12 total_params = sum(p.numel() for p in model.parameters())
13 print(f"Total parameters: {total_params}") # 11 (10 weights
14 + 1 bias)
15
16 # Access specific parameters
17 print(model.linear.weight) # The weight matrix
18 print(model.linear.bias) # The bias vector

```

3.3.3 Using nn.Sequential

For simple feed-forward architectures, `nn.Sequential` is convenient:

```

1 # Method 1: Pass modules as arguments
2 model = nn.Sequential(
3     nn.Linear(10, 50),
4     nn.ReLU(),
5     nn.Linear(50, 20),
6     nn.ReLU(),
7     nn.Linear(20, 1)
8 )

```

```
9
10 # Method 2: Use OrderedDict for named layers
11 from collections import OrderedDict
12
13 model = nn.Sequential(OrderedDict([
14     ('fc1', nn.Linear(10, 50)),
15     ('relu1', nn.ReLU()),
16     ('fc2', nn.Linear(50, 20)),
17     ('relu2', nn.ReLU()),
18     ('output', nn.Linear(20, 1))
19 ]))
20
21 # Use it
22 x = torch.randn(32, 10)
23 y = model(x)
24
25 # Access layers
26 print(model[0]) # First layer (fc1)
27 print(model.fc1) # Same, if using OrderedDict
```

When to use Sequential:

- Simple feed-forward architecture
- No branching or skip connections
- No custom logic in forward pass

When to use custom Module:

- Complex architectures (ResNets, U-Nets)
- Multiple inputs/outputs
- Custom forward logic (if statements, loops)
- Need to access intermediate values

3.3.4 Custom Modules with Multiple Layers

```

1 class MLP(nn.Module):
2     """Multi-Layer Perceptron with configurable hidden layers
3     ."""
4
5     def __init__(self, input_dim, hidden_dims, output_dim):
6         """
7         Args:
8             input_dim: Input feature dimension
9             hidden_dims: List of hidden layer dimensions
10            output_dim: Output dimension
11        """
12        super().__init__()
13
14        # Build layers dynamically
15        layers = []
16        prev_dim = input_dim
17
18        for hidden_dim in hidden_dims:
19            layers.append(nn.Linear(prev_dim, hidden_dim))
20            layers.append(nn.ReLU())
21            prev_dim = hidden_dim
22
23        # Output layer (no activation)
24        layers.append(nn.Linear(prev_dim, output_dim))
25
26        # Store as Sequential
27        self.network = nn.Sequential(*layers)
28
29    def forward(self, x):
30        return self.network(x)
31
32    # Create a 3-layer network
33    model = MLP(input_dim=10, hidden_dims=[50, 30], output_dim=1)
34
35    # The architecture is: 10 -> 50 -> 30 -> 1
36    # With ReLU after first two layers

```

3.3.5 Parameters vs Buffers

Parameters: Learnable tensors (updated by optimizer)

Buffers: Non-learnable tensors (not updated by optimizer, but saved with model)

```

1 class MyModule(nn.Module):
2     def __init__(self):
3         super().__init__()
4
5         # Learnable parameter
6         self.weight = nn.Parameter(torch.randn(10, 10))
7
8         # Non-learnable buffer (e.g., running statistics)
9         self.register_buffer('running_mean', torch.zeros(10))
10
11    def forward(self, x):
12        # Use both in computation
13        return x @ self.weight + self.running_mean
14
15    model = MyModule()
16
17    # Parameters (will be updated by optimizer)

```

```

18 print(len(list(model.parameters()))) # 1 (weight only)
19
20 # Buffers (won't be updated by optimizer)
21 print(len(list(model.buffers()))) # 1 (running_mean)
22
23 # Both are saved in state_dict
24 print(model.state_dict().keys())
25 # dict_keys(['weight', 'running_mean'])

```

When to use buffers:

- Running statistics in BatchNorm
- Fixed embeddings
- Constants that should be saved with the model
- Anything that needs to move with the model to GPU but shouldn't be trained

3.3.6 Training vs Evaluation Mode

Some layers behave differently during training and evaluation:

```

1 model = nn.Sequential(
2     nn.Linear(10, 50),
3     nn.ReLU(),
4     nn.Dropout(0.5), # Randomly drops 50% of values during
      training
5     nn.Linear(50, 1)
6 )
7
8 # Training mode (default)
9 model.train()
10 x = torch.randn(5, 10)
11 print(model(x)) # Dropout is active
12
13 # Evaluation mode
14 model.eval()
15 print(model(x)) # Dropout is disabled
16
17 # Check current mode
18 print(model.training) # False when in eval mode

```

[Always Set Mode Correctly] Forgetting to call `model.eval()` during inference can lead to incorrect results (dropout still active, batch norm using batch statistics instead of running statistics). Always:

```

1 model.train() # Before training
2 model.eval()  # Before evaluation/inference

```

3.3.7 Moving Models to GPU

```

1 device = torch.device('cuda' if torch.cuda.is_available()
2   else 'cpu')
3 # Create model
4 model = MLP(10, [50, 30], 1)
5
6 # Move to GPU
7 model = model.to(device)
8
9 # All parameters and buffers are now on GPU
10 print(next(model.parameters()).device) # cuda:0
11
12 # Input must also be on GPU
13 x = torch.randn(32, 10, device=device)
14 y = model(x) # Works! Both model and input on same device
15
16 # Moving back to CPU
17 model = model.to('cpu')

```

[Device Management Pattern] Common pattern for device-agnostic code:

```

1 device = torch.device('cuda' if torch.cuda.is_available()
2   else 'cpu')
3 model = model.to(device)
4
5 for batch in dataloader:
6     x, y = batch
7     x, y = x.to(device), y.to(device)
8     output = model(x)
9     # ... rest of training loop

```

3.3.8 Saving and Loading Models

```

1 # Save model weights
2 model = MLP(10, [50, 30], 1)
3 torch.save(model.state_dict(), 'model_weights.pth')
4
5 # Load weights into new model
6 new_model = MLP(10, [50, 30], 1) # Must have same
7   architecture!
8 new_model.load_state_dict(torch.load('model_weights.pth'))
9
10 # Save entire model (architecture + weights)
11 torch.save(model, 'entire_model.pth')
12 loaded_model = torch.load('entire_model.pth')
13
14 # For inference, set to eval mode
15 loaded_model.eval()

```

Best practice: Save `state_dict()` rather than entire model. This is more flexible and portable.

```

1 # Better: Save with additional info
2 checkpoint = {
3     'model_state_dict': model.state_dict(),
4     'optimizer_state_dict': optimizer.state_dict(),
5     'epoch': epoch,
6     'loss': loss,
7 }

```

```

8 torch.save(checkpoint, 'checkpoint.pth')
9
10 # Load
11 checkpoint = torch.load('checkpoint.pth')
12 model.load_state_dict(checkpoint['model_state_dict'])
13 optimizer.load_state_dict(checkpoint['optimizer_state_dict'])

```

3.3.9 Parameter Initialization

PyTorch layers have default initialization, but you often want custom initialization:

```

1 import torch.nn.init as init
2
3 class InitializedMLP(nn.Module):
4     def __init__(self, input_dim, hidden_dim, output_dim):
5         super().__init__()
6         self.fc1 = nn.Linear(input_dim, hidden_dim)
7         self.fc2 = nn.Linear(hidden_dim, output_dim)
8         self.relu = nn.ReLU()
9
10        # Custom initialization
11        self._initialize_weights()
12
13    def _initialize_weights(self):
14        # Xavier initialization for fc1
15        init.xavier_uniform_(self.fc1.weight)
16        init.zeros_(self.fc1.bias)
17
18        # He initialization for fc2 (better for ReLU)
19        init.kaiming_normal_(self.fc2.weight, nonlinearity='
relu')
20        init.zeros_(self.fc2.bias)
21
22    def forward(self, x):
23        x = self.relu(self.fc1(x))
24        x = self.fc2(x)
25        return x

```

Common initialization strategies:

Method	Use Case	Function
Xavier/Glorot	Sigmoid/Tanh	<code>init.xavier_uniform_</code>
He/Kaiming	ReLU	<code>init.kaiming_normal_</code>
Zeros	Biases	<code>init.zeros_</code>
Constant	Set to specific value	<code>init.constant_</code>
Normal	Custom std dev	<code>init.normal_</code>
Orthogonal	RNNs	<code>init.orthogonal_</code>

Table 2: Common weight initialization methods

[Why Initialization Matters] Poor initialization can cause:

- **Vanishing gradients:** Weights too small, gradients die out
- **Exploding gradients:** Weights too large, gradients explode
- **Slow convergence:** Network takes forever to learn
- **Dead neurons:** ReLUs output zero and never recover

Use Xavier for sigmoid/tanh activations, He for ReLU activations.

3.3.10 Module Hooks for Debugging

Hooks let you inspect intermediate values during forward/backward pass:

```

1 model = nn.Sequential(
2     nn.Linear(10, 50),
3     nn.ReLU(),
4     nn.Linear(50, 1)
5 )
6
7 # Forward hook: called after forward pass
8 def forward_hook(module, input, output):
9     print(f"Forward: {output.shape}")
10
11 # Register hook on second layer
12 handle = model[2].register_forward_hook(forward_hook)
13
14 # Forward pass
15 x = torch.randn(5, 10)
16 y = model(x) # Prints: Forward: torch.Size([5, 1])
17
18 # Remove hook when done
19 handle.remove()
20
21 # Backward hook: called during backward pass
22 def backward_hook(module, grad_input, grad_output):
23     print(f"Grad output shape: {grad_output[0].shape}")
24
25 handle = model[2].register_backward_hook(backward_hook)
26 y = model(x)
27 y.sum().backward() # Prints grad shape
28
29 handle.remove()

```

Use cases for hooks:

- Debugging shape mismatches
- Visualizing activations
- Gradient flow analysis
- Feature extraction from intermediate layers

3.4 Implementation: Practical Patterns

3.4.1 Freezing Layers

Sometimes you want to freeze part of the model (transfer learning):

```

1 # Freeze all parameters
2 for param in model.parameters():
3     param.requires_grad = False
4
5 # Freeze specific layer
6 for param in model.layer1.parameters():
7     param.requires_grad = False
8
9 # Only train the last layer
10 model = nn.Sequential(
11     nn.Linear(10, 50),
12     nn.ReLU(),
13     nn.Linear(50, 1)
14 )

```

```

15
16 # Freeze first two layers
17 for param in model[:2].parameters():
18     param.requires_grad = False
19
20 # Now only model[2].parameters() will be updated
21 optimizer = torch.optim.Adam(
22     filter(lambda p: p.requires_grad, model.parameters()),
23     lr=0.001
24 )

```

3.4.2 Multiple Inputs/Outputs

```

1 class MultiInputModel(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.branch1 = nn.Linear(10, 50)
5         self.branch2 = nn.Linear(5, 50)
6         self.combine = nn.Linear(100, 1)
7
8     def forward(self, x1, x2):
9         """
10        x1: shape (batch, 10)
11        x2: shape (batch, 5)
12        """
13        out1 = self.branch1(x1)
14        out2 = self.branch2(x2)
15        combined = torch.cat([out1, out2], dim=1) # (batch,
16        100)
17        return self.combine(combined)
18
19 model = MultiInputModel()
20 x1 = torch.randn(32, 10)
21 x2 = torch.randn(32, 5)
22 y = model(x1, x2)

```

3.4.3 Custom Layer with Parameters

```

1 class ScaledLinear(nn.Module):
2     """Linear layer with learnable scaling factor."""
3
4     def __init__(self, in_features, out_features):
5         super().__init__()
6         self.linear = nn.Linear(in_features, out_features)
7         # Custom learnable parameter
8         self.scale = nn.Parameter(torch.ones(1))
9
10    def forward(self, x):
11        return self.scale * self.linear(x)
12
13 layer = ScaledLinear(10, 5)
14
15 # Parameters include linear weights/bias AND scale
16 for name, param in layer.named_parameters():
17     print(f"{name}: {param.shape}")
18 # linear.weight: torch.Size([5, 10])
19 # linear.bias: torch.Size([5])
20 # scale: torch.Size([1])

```

3.5 Exercises

3.1: Simple Module - ★★

Goal: Create your first custom module.

Implement a `TwoLayerNet` module:

1. Two linear layers: $\text{input_dim} \rightarrow \text{hidden_dim} \rightarrow \text{output_dim}$
2. ReLU activation between them
3. No activation after the output layer
4. Test with $\text{input_dim}=20$, $\text{hidden_dim}=50$, $\text{output_dim}=10$
5. Verify the output shape for a batch of 16 samples

Starter code:

```

1 import torch
2 import torch.nn as nn
3
4 class TwoLayerNet(nn.Module):
5     def __init__(self, input_dim, hidden_dim, output_dim):
6         super().__init__()
7         # Your code here
8
9     def forward(self, x):
10        # Your code here
11        pass
12
13 # Test
14 model = TwoLayerNet(20, 50, 10)
15 x = torch.randn(16, 20)
16 y = model(x)
17 print(y.shape) # Should be torch.Size([16, 10])

```

3.2: Parameter Counting - ★★

Goal: Understand model size.

1. Create a model: $100 \rightarrow 500 \rightarrow 200 \rightarrow 50 \rightarrow 10$ (all linear, ReLU between)
2. Count total parameters programmatically
3. Calculate manually and verify:
 - Layer 1: $100 \times 500 + 500 = 50,500$
 - Layer 2: $500 \times 200 + 200 = 100,200$
 - ... (continue)
4. Print the number of parameters in each layer

Hint: Use `param.numel()` to count elements in each parameter.

3.3: Sequential vs Custom - ★★★

Goal: Compare two approaches.

Build the same network two ways:

1. Using `nn.Sequential`
2. Using a custom `nn.Module` with explicit layers
3. Both should have: $10 \rightarrow 50 \rightarrow 30 \rightarrow 1$ with ReLU activations
4. Load the same random seed before creating each

5. Verify they produce identical outputs for the same input

Challenge: Can you implement a method to copy weights from one to the other?

3.4: Custom Initialization - ★★☆☆

Goal: Implement custom weight initialization.

1. Create a 3-layer MLP
2. Initialize all weights with He initialization
3. Initialize all biases to zero
4. Create another identical model with default initialization
5. Train both on a simple task (e.g., fit $y = x^2$ for x in $[-1, 1]$)
6. Compare convergence speed (plot losses)

Hint: Use `nn.init.kaiming_normal_()` for He initialization.

3.5: Save and Load - ★★☆☆

Goal: Practice model persistence.

1. Create a model and train it for a few steps (any simple task)
2. Save the state dict
3. Create a new model instance
4. Load the weights
5. Verify both models produce identical outputs
6. Try saving/loading with `torch.save(model, ...)` and compare

3.6: Multi-Input Network - ★★☆☆

Goal: Build a more complex architecture.

Implement a model that takes two inputs and combines them:

1. Input 1: shape (batch, 20)
2. Input 2: shape (batch, 10)
3. Process each through separate 2-layer networks
4. Concatenate the results
5. Pass through a final layer to get output shape (batch, 1)
6. Test with random inputs

Extra challenge: Add skip connections (concatenate the original inputs with processed features before the final layer).

3.7: Frozen Layers - ★★☆☆

Goal: Practice transfer learning patterns.

1. Create a 4-layer network
2. "Pretrain" it on a simple task
3. Freeze the first two layers (`requires_grad=False`)
4. "Fine-tune" on a different task (only last two layers train)
5. Verify the first two layers' weights don't change
6. Compare with training all layers from scratch

Task idea: First task: classify points in quadrants, second task: classify by

distance from origin.

4 The Training Loop & Optimizers

4.1 Introduction: The Training Process

We now have all the pieces:

- **Tensors:** For data representation
- **Autograd:** For computing gradients
- **Modules:** For building models

Now we need to **train** the model—adjust its parameters to minimize a loss function. This involves:

1. Forward pass: Compute predictions
2. Loss computation: Measure how wrong we are
3. Backward pass: Compute gradients
4. Parameter update: Move in the direction that reduces loss
5. Repeat!

4.2 Theory: Gradient Descent and Optimization

4.2.1 The Optimization Problem

Training is an optimization problem. Given:

- Model f_θ with parameters θ
- Dataset $\{(x_i, y_i)\}_{i=1}^N$
- Loss function \mathcal{L}

We want to find:

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f_\theta(x_i), y_i)$$

4.2.2 Gradient Descent

The basic update rule:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}$$

where η is the **learning rate**.

Variants:

- **Batch Gradient Descent:** Use entire dataset (slow, memory intensive)
- **Stochastic Gradient Descent (SGD):** Use one sample (fast, noisy)
- **Mini-Batch SGD:** Use a batch of samples (sweet spot)

In practice, we always use mini-batch SGD (typically 32-256 samples per batch).

4.2.3 Learning Rate

The learning rate η is the most important hyperparameter:

- **Too large:** Training diverges (loss increases)
- **Too small:** Training is very slow
- **Just right:** Steady convergence

Typical ranges: 10^{-4} to 10^{-2} for Adam, 10^{-3} to 10^{-1} for SGD.

[Momentum and Adaptive Learning Rates] Modern optimizers go beyond basic gradient descent:

- **Momentum:** Accumulate gradients over time (helps escape local minima)
- **Adaptive rates:** Different learning rate for each parameter (Adam, RMSprop)

These techniques make training more stable and faster.

4.3 Implementation: The Standard Training Loop

4.3.1 Complete Training Example

Here's the canonical PyTorch training loop:

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4
5 # 1. CREATE MODEL
6 model = nn.Sequential(
7     nn.Linear(10, 50),
8     nn.ReLU(),
9     nn.Linear(50, 1)
10 )
11
12 # 2. DEFINE LOSS AND OPTIMIZER
13 criterion = nn.MSELoss()
14 optimizer = optim.Adam(model.parameters(), lr=0.001)
15
16 # 3. CREATE DUMMY DATA (replace with real DataLoader)
17 X_train = torch.randn(1000, 10)
18 y_train = torch.randn(1000, 1)
19
20 # 4. TRAINING LOOP
21 num_epochs = 100
22 batch_size = 32
23
24 for epoch in range(num_epochs):
25     # Shuffle data (simplified; use DataLoader in practice)
26     indices = torch.randperm(len(X_train))
27
28     epoch_loss = 0.0
29     for i in range(0, len(X_train), batch_size):
30         # Get batch
31         batch_indices = indices[i:i+batch_size]
32         X_batch = X_train[batch_indices]
33         y_batch = y_train[batch_indices]
34
35         # FORWARD PASS
36         predictions = model(X_batch)
37         loss = criterion(predictions, y_batch)
38
39         # BACKWARD PASS
40         optimizer.zero_grad() # Clear old gradients
41         loss.backward()        # Compute new gradients
42
43         # UPDATE PARAMETERS
44         optimizer.step()
45
46     epoch_loss += loss.item()

```

```
47
48     # Print progress
49     if (epoch + 1) % 10 == 0:
50         avg_loss = epoch_loss / (len(X_train) / batch_size)
51         print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {
            avg_loss:.4f}')
```

[The Five Steps of Training] Every training loop has these steps:

1. `optimizer.zero_grad()`: Clear old gradients
2. `output = model(input)`: Forward pass
3. `loss = criterion(output, target)`: Compute loss
4. `loss.backward()`: Compute gradients
5. `optimizer.step()`: Update parameters

Always in this order!

4.3.2 Loss Functions

Different tasks require different loss functions:

Regression (continuous output):

```

1 # Mean Squared Error (L2 loss)
2 criterion = nn.MSELoss()
3
4 # Mean Absolute Error (L1 loss, more robust to outliers)
5 criterion = nn.L1Loss()
6
7 # Smooth L1 (Huber loss, combines L1 and L2)
8 criterion = nn.SmoothL1Loss()

```

Binary Classification (0 or 1):

```

1 # Binary Cross Entropy (output must be in [0,1], use sigmoid)
2 criterion = nn.BCELoss()
3
4 # BCE with built-in sigmoid (more numerically stable)
5 criterion = nn.BCEWithLogitsLoss() # Preferred!
6
7 # Example
8 model = nn.Sequential(nn.Linear(10, 1)) # No sigmoid in
   model
9 criterion = nn.BCEWithLogitsLoss()
10 logits = model(x) # Raw outputs
11 loss = criterion(logits, targets) # Applies sigmoid
   internally

```

Multi-Class Classification (one of K classes):

```

1 # Cross Entropy Loss (combines LogSoftmax + NLLLoss)
2 criterion = nn.CrossEntropyLoss()
3
4 # Example
5 model = nn.Sequential(nn.Linear(10, 5)) # 5 classes, no
   softmax!
6 criterion = nn.CrossEntropyLoss()
7
8 logits = model(x) # Shape: (batch, 5), raw scores
9 targets = torch.tensor([0, 2, 1, 4, 3]) # Class indices
10
11 loss = criterion(logits, targets) # Applies softmax
   internally

```

[Don't Apply Softmax Before CrossEntropyLoss] `nn.CrossEntropyLoss` expects **raw logits**, not probabilities. It applies softmax internally for numerical stability. Applying softmax yourself will give wrong results!

WRONG:

```

1 output = torch.softmax(model(x), dim=1)
2 loss = nn.CrossEntropyLoss()(output, targets) # INCORRECT!

```

CORRECT:

```

1 output = model(x) # Raw logits
2 loss = nn.CrossEntropyLoss()(output, targets) # Correct!

```

Multi-Label Classification (multiple classes can be active):

```

1 # Use BCE with Logits for each class independently
2 criterion = nn.BCEWithLogitsLoss()
3
4 model = nn.Sequential(nn.Linear(10, 5)) # 5 binary outputs
5 logits = model(x) # Shape: (batch, 5)
6 targets = torch.tensor([[1, 0, 1, 0, 0], # Multiple 1s
7                        allowed
8                        [0, 1, 1, 1, 0]]) # Shape: (batch,
9                        5)
10 loss = criterion(logits, targets.float())

```

4.3.3 Optimizers

PyTorch provides many optimizers. Here are the most common:

Stochastic Gradient Descent (SGD):

```

1 # Basic SGD
2 optimizer = optim.SGD(model.parameters(), lr=0.01)
3
4 # SGD with momentum (recommended)
5 optimizer = optim.SGD(model.parameters(), lr=0.01, momentum
6                        =0.9)
7
8 # SGD with momentum and weight decay (L2 regularization)
9 optimizer = optim.SGD(model.parameters(), lr=0.01,
10                       momentum=0.9, weight_decay=1e-4)

```

Adam (Adaptive Moment Estimation):

```

1 # Adam: adaptive learning rates per parameter
2 optimizer = optim.Adam(model.parameters(), lr=0.001)
3
4 # Adam with weight decay
5 optimizer = optim.Adam(model.parameters(), lr=0.001,
6                       weight_decay=1e-4)

```

AdamW (Adam with decoupled weight decay):

```

1 # AdamW: better weight decay than Adam
2 optimizer = optim.AdamW(model.parameters(), lr=0.001,
3                          weight_decay=0.01)

```

[Which Optimizer to Use?] **For most tasks:** Start with Adam or AdamW with lr=0.001

Use Adam/AdamW when:

- You want fast initial progress
- You're prototyping
- You have many hyperparameters (adaptive rates help)

Use SGD with momentum when:

- You need the absolute best final performance
- You're training CNNs for vision tasks
- You can afford to tune learning rate carefully

General rule: Adam for fast development, SGD for final tuning.

4.3.4 Learning Rate Schedules

Instead of fixed learning rate, reduce it during training:

```

1 import torch.optim.lr_scheduler as lr_scheduler
2
3 optimizer = optim.Adam(model.parameters(), lr=0.001)
4
5 # Step decay: reduce LR every step_size epochs
6 scheduler = lr_scheduler.StepLR(optimizer, step_size=30,
7                                 gamma=0.1)
8
9 # Exponential decay
10 scheduler = lr_scheduler.ExponentialLR(optimizer, gamma=0.95)
11
12 # Reduce on plateau: reduce when loss stops improving
13 scheduler = lr_scheduler.ReduceLROnPlateau(optimizer, mode='
14 min',
15                                             factor=0.5,
16                                             patience=10)
17
18 # Cosine annealing: smooth decay to zero
19 scheduler = lr_scheduler.CosineAnnealingLR(optimizer, T_max
20 =100)
21
22 # Usage in training loop
23 for epoch in range(num_epochs):
24     train_one_epoch()
25
26     # Update learning rate
27     scheduler.step() # For most schedulers
28
29     # For ReduceLROnPlateau, pass validation loss
30     # scheduler.step(val_loss)
31
32     # Check current LR
33     current_lr = optimizer.param_groups[0]['lr']
34     print(f"Current LR: {current_lr}")

```

When to use schedules:

- Long training runs (>100 epochs)
- When loss plateaus
- Fine-tuning pretrained models

When not to use:

- Short training (<50 epochs)
- Already using Adam with good convergence
- Adds complexity without clear benefit

4.3.5 Gradient Clipping

Prevent exploding gradients by limiting their magnitude:

```

1 # Training loop with gradient clipping
2 for epoch in range(num_epochs):
3     for batch in dataloader:
4         optimizer.zero_grad()
5
6         output = model(batch)

```



```
7         loss = criterion(output, targets)
8         loss.backward()
9
10        # Clip gradients by norm
11        torch.nn.utils.clip_grad_norm_(model.parameters(),
max_norm=1.0)
12
13        optimizer.step()
```

Two methods:

```
1 # Clip by norm (preferred): scale down if norm > max_norm
2 torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm
=1.0)
3
4 # Clip by value: clamp each gradient to [-clip_value,
clip_value]
5 torch.nn.utils.clip_grad_value_(model.parameters(),
clip_value=0.5)
```

[When to Use Gradient Clipping] **Always use for:**

- RNNs and LSTMs (prone to exploding gradients)
- Transformers
- Any sequence model

Often useful for:

- Deep networks (>50 layers)
- Training instability

Typical values: max_norm=1.0 to 5.0

4.3.6 Validation Loop

Always evaluate on held-out data:

```

1 def train_epoch(model, dataloader, criterion, optimizer,
2   device):
3     model.train() # Set to training mode
4     total_loss = 0.0
5
6     for batch_idx, (data, target) in enumerate(dataloader):
7         data, target = data.to(device), target.to(device)
8
9         optimizer.zero_grad()
10        output = model(data)
11        loss = criterion(output, target)
12        loss.backward()
13        optimizer.step()
14
15        total_loss += loss.item()
16
17    return total_loss / len(dataloader)
18
19 def validate(model, dataloader, criterion, device):
20     model.eval() # Set to evaluation mode
21     total_loss = 0.0
22     correct = 0
23     total = 0
24
25     with torch.no_grad(): # Disable gradient computation
26         for data, target in dataloader:
27             data, target = data.to(device), target.to(device)
28
29             output = model(data)
30             loss = criterion(output, target)
31             total_loss += loss.item()
32
33             # For classification: compute accuracy
34             _, predicted = torch.max(output, 1)
35             total += target.size(0)
36             correct += (predicted == target).sum().item()
37
38     avg_loss = total_loss / len(dataloader)
39     accuracy = 100 * correct / total
40     return avg_loss, accuracy
41
42 # Main training loop
43 for epoch in range(num_epochs):
44     train_loss = train_epoch(model, train_loader, criterion,
45                             optimizer, device)
46     val_loss, val_acc = validate(model, val_loader, criterion,
47                                device)
48
49     print(f'Epoch {epoch+1}: Train Loss={train_loss:.4f}, '
50           f'Val Loss={val_loss:.4f}, Val Acc={val_acc:.2f}%')
```

[Don't Forget `model.eval()` and `torch.no_grad()`] When validating:

- `model.eval()`: Disables dropout, uses running stats for batch norm
- `torch.no_grad()`: Saves memory and computation

Forgetting these can lead to:

- Incorrect validation metrics
- Out of memory errors
- Slower evaluation

4.3.7 Early Stopping

Stop training when validation loss stops improving:

```

1 class EarlyStopping:
2     def __init__(self, patience=10, min_delta=0):
3         """
4         patience: how many epochs to wait after last
5         improvement
6         min_delta: minimum change to qualify as improvement
7         """
8         self.patience = patience
9         self.min_delta = min_delta
10        self.counter = 0
11        self.best_loss = None
12        self.early_stop = False
13
14    def __call__(self, val_loss):
15        if self.best_loss is None:
16            self.best_loss = val_loss
17        elif val_loss > self.best_loss - self.min_delta:
18            self.counter += 1
19            if self.counter >= self.patience:
20                self.early_stop = True
21        else:
22            self.best_loss = val_loss
23            self.counter = 0
24
25    # Usage
26    early_stopping = EarlyStopping(patience=10)
27
28    for epoch in range(num_epochs):
29        train_loss = train_epoch(...)
30        val_loss, val_acc = validate(...)
31
32        early_stopping(val_loss)
33        if early_stopping.early_stop:
34            print(f"Early stopping at epoch {epoch+1}")
35            break

```

4.4 Implementation: Debugging Training

4.4.1 When Training Goes Wrong

Problem: Loss is NaN

```

1 # Check for NaN
2 if torch.isnan(loss):
3     print("NaN detected!")
4     print(f"Input contains NaN: {torch.isnan(data).any()}")
5     print(f"Output contains NaN: {torch.isnan(output).any()}")
6
7 # Check gradients
8 for name, param in model.named_parameters():
9     if param.grad is not None:
10        print(f"{name} grad: {param.grad.abs().max()}")

```

Common causes:

- Learning rate too high
- Numerical instability (use BCEWithLogitsLoss, not BCELoss)
- Division by zero
- Log of zero or negative number

Problem: Loss doesn't decrease

```

1 # Sanity checks
2 print("Model output range:", output.min().item(), output.max()
3     .item())
4 print("Target range:", target.min().item(), target.max().item()
5     .item())
6
7 # Check if gradients are flowing
8 for name, param in model.named_parameters():
9     if param.grad is not None:
10        print(f"{name}: grad mean = {param.grad.abs().mean()
11            :.6f}")

```

Common causes:

- Learning rate too low
- Wrong loss function
- Forgot `optimizer.zero_grad()`
- Model architecture issue (dead ReLUs, wrong dimensions)
- Data not normalized

Problem: Loss explodes

- Learning rate too high (reduce by 10x)
- Use gradient clipping
- Check for large batch effects (try smaller batches)

Problem: Overfitting (train loss « val loss)

```

1 # Add regularization
2 model = nn.Sequential(
3     nn.Linear(10, 50),
4     nn.ReLU(),

```

```

5     nn.Dropout(0.5), # Add dropout
6     nn.Linear(50, 1)
7 )
8
9 # Add weight decay
10 optimizer = optim.Adam(model.parameters(), lr=0.001,
    weight_decay=1e-4)
11
12 # Use data augmentation (for images)
13 # Reduce model size
14 # Get more data

```

Problem: Underfitting (both losses high)

- Increase model capacity (more layers, more neurons)
- Train longer
- Increase learning rate
- Remove regularization
- Check if task is actually learnable

4.4.2 Monitoring Training

```

1 # Track metrics
2 train_losses = []
3 val_losses = []
4
5 for epoch in range(num_epochs):
6     train_loss = train_epoch(...)
7     val_loss, val_acc = validate(...)
8
9     train_losses.append(train_loss)
10    val_losses.append(val_loss)
11
12    # Plot periodically
13    if epoch % 10 == 0:
14        import matplotlib.pyplot as plt
15        plt.plot(train_losses, label='Train')
16        plt.plot(val_losses, label='Val')
17        plt.xlabel('Epoch')
18        plt.ylabel('Loss')
19        plt.legend()
20        plt.show()

```

4.5 Exercises

4.1: Basic Training Loop - ★★

Goal: Implement a complete training loop from scratch.

Task: Fit a neural network to $y = \sin(x)$ for $x \in [0, 2\pi]$

1. Generate 1000 training samples: x uniform in $[0, 2\pi]$, $y = \sin(x) + \text{small noise}$
2. Create a 3-layer MLP: $1 \rightarrow 50 \rightarrow 50 \rightarrow 1$ with ReLU
3. Use MSE loss and Adam optimizer
4. Train for 100 epochs, batch size 32
5. Print loss every 10 epochs
6. Plot predictions vs true function

Starter code:

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import numpy as np
5
6 # Generate data
7 x = torch.linspace(0, 2*np.pi, 1000).reshape(-1, 1)
8 y = torch.sin(x) + 0.1 * torch.randn_like(x)
9
10 # Your code here

```

4.2: Loss Function Comparison - ★★★

Goal: Understand different loss functions.

Train the same model with different losses:

1. MSE Loss (L2)
2. L1 Loss
3. Smooth L1 Loss

Add outliers to the data (10% of points with large noise) and compare:

- Final training loss
- Robustness to outliers
- Convergence speed

Which loss is most robust to outliers?

4.3: Optimizer Comparison - ★★★

Goal: Compare optimizer behavior.

Train the same model with:

1. SGD (lr=0.01)
2. SGD with momentum (lr=0.01, momentum=0.9)
3. Adam (lr=0.001)
4. AdamW (lr=0.001)

Plot all loss curves on the same graph. Observations:

- Which converges fastest initially?
- Which achieves lowest final loss?
- How does momentum help SGD?

4.4: Learning Rate Experiments - ★★★**Goal:** Understand learning rate impact.

Train the same model with different learning rates:

- 10^{-4} , 10^{-3} , 10^{-2} , 10^{-1} , 1.0

For each:

1. Track loss over epochs
2. Identify: too small, just right, too large
3. Observe divergence (NaN loss) for very large LR

Extra: Implement a learning rate finder (sweep from small to large, stop when loss diverges).**4.5: Overfitting Detection - ★★★★★****Goal:** Identify and prevent overfitting.

1. Generate a small dataset (100 samples)
2. Split into train (80) and validation (20)
3. Train a large model (can easily overfit)
4. Track both train and validation loss
5. Identify when overfitting starts (losses diverge)
6. Add regularization:
 - Dropout
 - Weight decay
 - Early stopping
7. Compare final validation performance

4.6: Complete Training Pipeline - ★★★★★**Goal:** Build a production-ready training pipeline.

Implement a complete pipeline with:

1. Data splitting (train/val/test)
2. Training loop with progress bars
3. Validation after each epoch
4. Model checkpointing (save best model)
5. Early stopping
6. Learning rate scheduling
7. Tensorboard logging (optional)
8. Final test set evaluation

Test on a real dataset (e.g., sklearn's `make_classification`).**Hint:** Use `tqdm` for progress bars:

```

1 from tqdm import tqdm
2
3 for epoch in tqdm(range(num_epochs)):
4     # training code

```

4.7: Debugging Challenge - ★★★★★

Goal: Practice debugging common training issues.

Given buggy training code (deliberately broken):

```
1 # This code has 5 bugs - find and fix them!
2 model = nn.Sequential(nn.Linear(10, 1), nn.Softmax())
3 criterion = nn.CrossEntropyLoss()
4 optimizer = optim.Adam(model.parameters(), lr=10.0)
5
6 for epoch in range(100):
7     output = model(X_train)
8     loss = criterion(output, y_train)
9     loss.backward()
10    optimizer.step()
11    print(f"Epoch {epoch}, Loss: {loss.item()}")
```

Find and fix:

1. Wrong activation before CrossEntropyLoss
2. Missing `optimizer.zero_grad()`
3. Learning rate too high
4. No train/eval modes
5. No validation set

5 Data Loading & Preprocessing

5.1 Introduction: Why Proper Data Loading Matters

So far, we've used simple tensors for data. In practice, you'll have:

- Large datasets that don't fit in memory
- Complex preprocessing pipelines
- Data on disk (files, databases)
- Need for efficient batching and shuffling
- Different data types (images, time series, point clouds, meshes)

PyTorch provides `Dataset` and `DataLoader` to handle this elegantly:

- **Dataset:** Defines how to access individual samples
- **DataLoader:** Handles batching, shuffling, parallel loading

5.2 Theory: Dataset and DataLoader Architecture

[The Data Pipeline]

Raw Data → Dataset → DataLoader → Model
(access) (batch, shuffle, parallel)

Dataset answers: "How do I get sample i ?"

DataLoader answers: "How do I create batches efficiently?"

Dataset Requirements:

A valid dataset must implement:

1. `__len__()`: Return total number of samples
2. `__getitem__(idx)`: Return sample at index `idx`

DataLoader Benefits:

- Automatic batching
- Shuffling for training
- Parallel data loading (multiple workers)
- Automatic memory pinning for GPU transfer
- Handling variable-length sequences

5.3 Implementation: Basic Usage

5.3.1 Using Built-in Datasets

PyTorch provides common datasets in `torchvision.datasets`:

```
1 import torch
2 from torch.utils.data import DataLoader
3 import torchvision
4 import torchvision.transforms as transforms
5
6 # Load MNIST dataset
7 transform = transforms.ToTensor() # Convert PIL Image to
   tensor
8
9 train_dataset = torchvision.datasets.MNIST(
```

```

10     root='./data',
11     train=True,
12     download=True,
13     transform=transform
14 )
15
16 test_dataset = torchvision.datasets.MNIST(
17     root='./data',
18     train=False,
19     download=True,
20     transform=transform
21 )
22
23 # Create DataLoaders
24 train_loader = DataLoader(
25     train_dataset,
26     batch_size=64,
27     shuffle=True,          # Shuffle training data
28     num_workers=2,        # Parallel loading
29     pin_memory=True       # Faster GPU transfer
30 )
31
32 test_loader = DataLoader(
33     test_dataset,
34     batch_size=64,
35     shuffle=False,        # Don't shuffle test data
36     num_workers=2
37 )
38
39 # Iterate through batches
40 for batch_idx, (images, labels) in enumerate(train_loader):
41     # images: (64, 1, 28, 28) - batch of 64 grayscale 28x28
42     # labels: (64,) - batch of 64 labels
43     print(f"Batch {batch_idx}: {images.shape}, {labels.shape}")
44     break

```

[DataLoader Parameters] Key parameters:

- **batch_size**: Number of samples per batch (32-256 typically)
- **shuffle**: Randomize order (True for training, False for evaluation)
- **num_workers**: Parallel processes for loading (0=main process, 2-4 for speed)
- **pin_memory**: Faster CPU→GPU transfer (True if using GPU)
- **drop_last**: Drop incomplete last batch (useful for batch norm)

5.3.2 Creating Custom Datasets

For your own data, create a custom Dataset:

```

1 from torch.utils.data import Dataset
2
3 class SimpleDataset(Dataset):
4     """Example: Dataset from numpy arrays."""
5
6     def __init__(self, X, y):
7         """
8         Args:
9             X: numpy array or tensor of features
10            y: numpy array or tensor of labels
11        """
12        self.X = torch.FloatTensor(X) if not isinstance(X,
13 torch.Tensor) else X
14        self.y = torch.FloatTensor(y) if not isinstance(y,
15 torch.Tensor) else y
16
17    def __len__(self):
18        """Return total number of samples."""
19        return len(self.X)
20
21    def __getitem__(self, idx):
22        """Return sample at index idx."""
23        return self.X[idx], self.y[idx]
24
25    # Usage
26    import numpy as np
27
28    X = np.random.randn(1000, 10)
29    y = np.random.randn(1000, 1)
30
31    dataset = SimpleDataset(X, y)
32    loader = DataLoader(dataset, batch_size=32, shuffle=True)
33
34    for batch_X, batch_y in loader:
35        print(batch_X.shape, batch_y.shape)
36        break

```

5.3.3 Dataset from Files

More realistic: loading data from disk:

```

1 import os
2 import pandas as pd
3 from PIL import Image
4
5 class ImageDataset(Dataset):
6     """Load images from a directory with CSV labels."""
7
8     def __init__(self, csv_file, img_dir, transform=None):
9         """
10        Args:
11            csv_file: Path to CSV with columns [filename,
12 label]
13            img_dir: Directory with images
14            transform: Optional transform to apply
15        """
16        self.data = pd.read_csv(csv_file)
17        self.img_dir = img_dir

```

```

17     self.transform = transform
18
19     def __len__(self):
20         return len(self.data)
21
22     def __getitem__(self, idx):
23         # Get image filename and label
24         img_name = self.data.iloc[idx, 0]
25         label = self.data.iloc[idx, 1]
26
27         # Load image
28         img_path = os.path.join(self.img_dir, img_name)
29         image = Image.open(img_path).convert('RGB')
30
31         # Apply transforms
32         if self.transform:
33             image = self.transform(image)
34
35         return image, label
36
37 # Usage with transforms
38 transform = transforms.Compose([
39     transforms.Resize((224, 224)),
40     transforms.ToTensor(),
41     transforms.Normalize(mean=[0.485, 0.456, 0.406],
42                           std=[0.229, 0.224, 0.225])
43 ])
44
45 dataset = ImageDataset('labels.csv', 'images/', transform=
46                        transform)
47 loader = DataLoader(dataset, batch_size=32, shuffle=True)

```

5.3.4 Dataset for Time Series

For scientific computing, time series are common:

```

1 class TimeSeriesDataset(Dataset):
2     """
3     Create windows from a long time series.
4     Useful for forecasting, where we predict future from past
5     .
6     """
7     def __init__(self, data, window_size, forecast_horizon):
8         """
9         Args:
10             data: 1D array of time series values
11             window_size: Length of input sequence
12             forecast_horizon: How many steps ahead to predict
13         """
14         self.data = torch.FloatTensor(data)
15         self.window_size = window_size
16         self.forecast_horizon = forecast_horizon
17
18     def __len__(self):
19         # Number of valid windows
20         return len(self.data) - self.window_size - self.
21         forecast_horizon + 1
22
23     def __getitem__(self, idx):
24         # Input: window_size points starting at idx

```

```
24         x = self.data[idx : idx + self.window_size]
25
26         # Target: forecast_horizon points after the input
27         y = self.data[idx + self.window_size :
28                       idx + self.window_size + self.
29               forecast_horizon]
30
31         return x, y
32
33 # Example: predict next 10 steps from past 50
34 data = np.sin(np.linspace(0, 100, 1000))
35 dataset = TimeSeriesDataset(data, window_size=50,
36                             forecast_horizon=10)
37
38 print(f"Dataset size: {len(dataset)}")
39 x, y = dataset[0]
40 print(f"Input shape: {x.shape}, Target shape: {y.shape}")
41 # Input shape: torch.Size([50]), Target shape: torch.Size
42 ([10])
43
44 loader = DataLoader(dataset, batch_size=32, shuffle=False)
45 # Note: Don't shuffle time series if order matters!
```

5.3.5 Dataset for Point Clouds

For 3D data (molecular structures, point clouds):

```

1 class PointCloudDataset(Dataset):
2     """Dataset for 3D point clouds."""
3
4     def __init__(self, data_dir, num_points=1024):
5         """
6         Args:
7             data_dir: Directory with .npy files, each
8             containing Nx3 points
9             num_points: Sample this many points from each
10            cloud
11            """
12            self.files = [os.path.join(data_dir, f)
13                           for f in os.listdir(data_dir) if f.
14                           endsuffix('.npy')]
15            self.num_points = num_points
16
17        def __len__(self):
18            return len(self.files)
19
20        def __getitem__(self, idx):
21            # Load point cloud (N x 3)
22            points = np.load(self.files[idx])
23
24            # Sample fixed number of points
25            if len(points) > self.num_points:
26                # Random sampling
27                indices = np.random.choice(len(points), self.
28                num_points,
29                replace=False)
30                points = points[indices]
31            elif len(points) < self.num_points:
32                # Pad with zeros
33                padding = np.zeros((self.num_points - len(points)
34                , 3))
35                points = np.vstack([points, padding])
36
37            return torch.FloatTensor(points)
38
39        # Usage
40        # dataset = PointCloudDataset('point_clouds/')
41        # loader = DataLoader(dataset, batch_size=8, shuffle=True)
42        # for batch in loader:
43        #     print(batch.shape) # (8, 1024, 3)

```

5.3.6 Data Transforms

Transforms process data on-the-fly:

```

1 from torchvision import transforms
2
3 # Common image transforms
4 transform = transforms.Compose([
5     # Resize to fixed size
6     transforms.Resize((224, 224)),
7
8     # Data augmentation (for training)
9     transforms.RandomHorizontalFlip(p=0.5),
10    transforms.RandomRotation(degrees=15),

```

```

11     transforms.ColorJitter(brightness=0.2, contrast=0.2),
12
13     # Convert to tensor
14     transforms.ToTensor(), # Converts to [0, 1] and (C, H, W
15     )
16
17     # Normalize (important!)
18     transforms.Normalize(mean=[0.485, 0.456, 0.406], #
19     ImageNet stats
20                             std=[0.229, 0.224, 0.225])
21 ]))
22
23 # For evaluation/test: no data augmentation
24 test_transform = transforms.Compose([
25     transforms.Resize((224, 224)),
26     transforms.ToTensor(),
27     transforms.Normalize(mean=[0.485, 0.456, 0.406],
28                             std=[0.229, 0.224, 0.225])
29 ])
30
31 # Custom transform as a callable
32 class AddGaussianNoise:
33     """Add Gaussian noise to tensor."""
34     def __init__(self, mean=0., std=0.1):
35         self.mean = mean
36         self.std = std
37
38     def __call__(self, tensor):
39         return tensor + torch.randn(tensor.size()) * self.std
40         + self.mean
41
42 # Use in Compose
43 transform = transforms.Compose([
44     transforms.ToTensor(),
45     AddGaussianNoise(std=0.05)
46 ])

```

[When to Apply Transforms] Transforms can be applied in two places:

In Dataset (___getitem___):

- Pro: Each sample transformed on-the-fly
- Pro: Can use data augmentation
- Con: Adds computation during training

Pre-computed and saved:

- Pro: Faster training (no transform overhead)
- Con: More disk space
- Con: Can't use random augmentation

For training with augmentation: use Dataset transforms. For inference: consider pre-computing.

5.3.7 Normalization Strategies

Normalization is crucial for training stability:

```

1  # Method 1: Standardization (zero mean, unit variance)
2  def compute_mean_std(dataset):
3      """Compute mean and std for normalization."""
4      loader = DataLoader(dataset, batch_size=64, shuffle=False)
5
6      mean = 0.0
7      std = 0.0
8      total_samples = 0
9
10     for data, _ in loader:
11         batch_samples = data.size(0)
12         data = data.view(batch_samples, data.size(1), -1)
13         mean += data.mean(2).sum(0)
14         std += data.std(2).sum(0)
15         total_samples += batch_samples
16
17     mean /= total_samples
18     std /= total_samples
19
20     return mean, std
21
22 # Usage
23 # mean, std = compute_mean_std(train_dataset)
24 # transform = transforms.Normalize(mean=mean, std=std)
25
26 # Method 2: Min-Max scaling to [0, 1] or [-1, 1]
27 class MinMaxScaler:
28     def __init__(self, data):
29         self.min = data.min()
30         self.max = data.max()
31
32     def transform(self, data):
33         return (data - self.min) / (self.max - self.min)
34
35     def inverse_transform(self, data):
36         return data * (self.max - self.min) + self.min
37
38 # Method 3: Per-feature normalization (for tabular data)
39 class StandardScaler:
40     def __init__(self, data):
41         """
42         data: (N, D) tensor
43         """
44         self.mean = data.mean(dim=0)
45         self.std = data.std(dim=0)
46
47     def transform(self, data):
48         return (data - self.mean) / (self.std + 1e-8)
49
50     def inverse_transform(self, data):
51         return data * self.std + self.mean
52
53 # Usage
54 X_train = torch.randn(1000, 10)
55 scaler = StandardScaler(X_train)
56 X_train_normalized = scaler.transform(X_train)
57 X_test_normalized = scaler.transform(X_test) # Use training

```


statistics!

[Always Use Training Statistics] When normalizing:

1. Compute mean/std from **training data only**
2. Apply the same transformation to validation and test data
3. **Never** compute statistics from test data!

Why? Using test statistics is a form of data leakage and gives overly optimistic results.

5.3.8 Handling Variable-Length Sequences

For sequences of different lengths (text, time series with variable duration):

```

1 from torch.nn.utils.rnn import pad_sequence,
  pack_padded_sequence
2
3 def collate_fn(batch):
4     """
5     Custom collate function for variable-length sequences.
6
7     batch: List of (sequence, label) tuples
8     sequences: List of tensors with different lengths
9     """
10    sequences, labels = zip(*batch)
11
12    # Get lengths
13    lengths = torch.LongTensor([len(seq) for seq in sequences
14    ])
15
16    # Pad sequences to same length
17    sequences_padded = pad_sequence(sequences, batch_first=
18    True,
19                                padding_value=0)
20
21    labels = torch.LongTensor(labels)
22
23    return sequences_padded, labels, lengths
24
25 # Example dataset with variable lengths
26 class VariableLengthDataset(Dataset):
27     def __init__(self, num_samples=100):
28         self.data = []
29         for _ in range(num_samples):
30             # Random length between 10 and 50
31             length = torch.randint(10, 50, (1,)).item()
32             seq = torch.randn(length, 5) # 5 features
33             label = torch.randint(0, 2, (1,)).item()
34             self.data.append((seq, label))
35
36     def __len__(self):
37         return len(self.data)
38
39     def __getitem__(self, idx):
40         return self.data[idx]
41
42 # Use with custom collate_fn
43 dataset = VariableLengthDataset()
44 loader = DataLoader(dataset, batch_size=8, collate_fn=
45                     collate_fn)

```

```
44 for sequences, labels, lengths in loader:
45     print(f"Padded sequences: {sequences.shape}") # (8,
      max_len, 5)
46     print(f"Lengths: {lengths}") # (8,)
47     break
```

5.3.9 Train/Validation/Test Split

Proper data splitting:

```

1 from torch.utils.data import random_split
2
3 # Method 1: random_split
4 dataset = SimpleDataset(X, y)
5 train_size = int(0.7 * len(dataset))
6 val_size = int(0.15 * len(dataset))
7 test_size = len(dataset) - train_size - val_size
8
9 train_dataset, val_dataset, test_dataset = random_split(
10     dataset, [train_size, val_size, test_size]
11 )
12
13 # Method 2: Using indices (more control)
14 from torch.utils.data import Subset
15 import numpy as np
16
17 indices = np.arange(len(dataset))
18 np.random.shuffle(indices)
19
20 train_idx = indices[:train_size]
21 val_idx = indices[train_size:train_size+val_size]
22 test_idx = indices[train_size+val_size:]
23
24 train_dataset = Subset(dataset, train_idx)
25 val_dataset = Subset(dataset, val_idx)
26 test_dataset = Subset(dataset, test_idx)
27
28 # Method 3: sklearn for stratified split (classification)
29 from sklearn.model_selection import train_test_split
30
31 X = np.random.randn(1000, 10)
32 y = np.random.randint(0, 3, 1000) # 3 classes
33
34 X_train, X_temp, y_train, y_temp = train_test_split(
35     X, y, test_size=0.3, stratify=y, random_state=42
36 )
37 X_val, X_test, y_val, y_test = train_test_split(
38     X_temp, y_temp, test_size=0.5, stratify=y_temp,
39     random_state=42
40 )
41
42 train_dataset = SimpleDataset(X_train, y_train)
43 val_dataset = SimpleDataset(X_val, y_val)
44 test_dataset = SimpleDataset(X_test, y_test)

```

5.3.10 Efficient Data Loading

For large datasets that don't fit in RAM:

```

1 class LazyLoadDataset(Dataset):
2     """Load data on-demand, don't store in memory."""
3
4     def __init__(self, file_list):
5         """
6         file_list: List of file paths
7         """
8         self.file_list = file_list
9

```

```

10     def __len__(self):
11         return len(self.file_list)
12
13     def __getitem__(self, idx):
14         # Load from disk only when needed
15         data = np.load(self.file_list[idx])
16         return torch.FloatTensor(data)
17
18 # For even larger datasets: memory-mapped arrays
19 class MemMapDataset(Dataset):
20     """Use memory-mapped arrays for huge datasets."""
21
22     def __init__(self, memmap_file, shape, dtype):
23         """
24         memmap_file: Path to .npy file
25         shape: Shape of the data array
26         dtype: Data type
27         """
28         self.data = np.memmap(memmap_file, dtype=dtype,
29                               mode='r', shape=shape)
30
31     def __len__(self):
32         return self.data.shape[0]
33
34     def __getitem__(self, idx):
35         # Read only the needed slice from disk
36         return torch.FloatTensor(self.data[idx])

```

[Optimizing DataLoader Performance] If data loading is slow:

1. **Increase num_workers:** 2-4 usually optimal
2. **Use pin_memory=True:** If using GPU
3. **Prefetch:** DataLoader automatically prefetches next batch
4. **Reduce transforms:** Pre-compute expensive transforms
5. **Use SSD:** Faster disk I/O
6. **Check CPU bottleneck:** If GPU is idle, data loading is the issue

Warning: Too many workers can hurt performance (overhead).

5.4 Exercises

5.1: Simple Dataset - ★★

Goal: Create your first custom dataset.

1. Create a synthetic dataset: $y = 2x_1 + 3x_2 + \epsilon$
2. Generate 1000 samples with $x_1, x_2 \sim \mathcal{N}(0, 1)$ and noise $\epsilon \sim \mathcal{N}(0, 0.1)$
3. Implement a custom `Dataset` class
4. Create a `DataLoader` with `batch_size=32`
5. Iterate through one epoch and print batch shapes

5.2: Train/Val/Test Split - ★★

Goal: Practice proper data splitting.

1. Generate a classification dataset (1000 samples, 3 classes)
2. Split into train (70%), validation (15%), test (15%)
3. Ensure the split is stratified (balanced classes in each split)
4. Create separate `DataLoaders` for each split
5. Verify class distributions are similar

Hint: Use `sklearn.model_selection.train_test_split` with `stratify`.

5.3: Time Series Dataset - ★★★

Goal: Handle sequential data.

1. Generate a noisy sine wave with 10,000 points
2. Create a dataset that returns windows of size 50 to predict the next 10 points
3. Implement the `TimeSeriesDataset` class
4. Train a simple MLP to predict future values
5. Visualize predictions vs ground truth

Challenge: Try with a real time series (stock prices, weather data).

5.4: Data Normalization - ★★★

Goal: Understand normalization impact.

1. Create a dataset with features on different scales:
 - Feature 1: range [0, 1]
 - Feature 2: range [0, 1000]
 - Feature 3: range [-100, 100]
2. Train a model **without** normalization
3. Train a model **with** standardization
4. Compare:
 - Training speed (epochs to converge)
 - Final loss
 - Gradient magnitudes

5.5: Data Augmentation - ★★★

Goal: Implement custom transforms.

1. Create a simple image dataset (or use MNIST)
2. Implement custom transforms:
 - Add Gaussian noise
 - Random scaling
 - Random shifts
3. Train models with and without augmentation
4. Compare validation performance
5. Visualize augmented samples

5.6: Variable-Length Sequences - ★★★★★

Goal: Handle sequences of different lengths.

1. Create a dataset of sequences with random lengths (10-100)
2. Implement a custom `collate_fn` to pad sequences
3. Create a `DataLoader` using this collate function
4. Train an RNN that handles the padded sequences
5. Use `pack_padded_sequence` for efficiency (we'll learn this properly in the RNN section)

Hint: The collate function should return sequences, lengths, and labels.

5.7: Complete Data Pipeline - ★★★★★

Goal: Build a production-ready data pipeline.

Create a complete pipeline with:

1. Custom dataset loading from CSV files
2. Train/val/test split with stratification
3. Different transforms for train (with augmentation) and test
4. Normalization using training statistics
5. Efficient `DataLoaders` with multiple workers
6. Save and load normalization parameters
7. Handle class imbalance with weighted sampling (optional)

Test on a real dataset (e.g., UCI datasets, Kaggle).

Part II

Deep Learning Architectures

6 Fully Connected Networks (MLPs)

6.1 Introduction: The Foundation of Deep Learning

Multi-Layer Perceptrons (MLPs), also called fully connected networks or feedforward networks, are the simplest and most fundamental deep learning architecture. Every neuron in one layer connects to every neuron in the next layer—hence "fully connected."

Despite their simplicity, MLPs are:

- **Universal approximators:** Can approximate any continuous function
- **Building blocks:** Components of more complex architectures
- **Highly effective:** For tabular data, function approximation, and many scientific computing tasks

When to use MLPs:

- Tabular data (features don't have spatial/temporal structure)
- Function approximation in scientific computing
- As components in larger architectures
- When you have relatively small input dimensions (<1000s of features)

When NOT to use MLPs:

- Images (use CNNs—exploit spatial structure)
- Sequences (use RNNs/Transformers—exploit temporal structure)
- Very high-dimensional inputs (too many parameters)

6.2 Theory: Understanding MLPs

6.2.1 Architecture

An MLP consists of:

- **Input layer:** Receives features $\mathbf{x} \in \mathbb{R}^{d_{in}}$
- **Hidden layers:** Transform representations
- **Output layer:** Produces predictions $\hat{\mathbf{y}} \in \mathbb{R}^{d_{out}}$

Single hidden layer MLP:

$$\mathbf{h} = \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\hat{\mathbf{y}} = \mathbf{W}_2 \mathbf{h} + \mathbf{b}_2$$

where:

- $\mathbf{W}_1 \in \mathbb{R}^{h \times d_{in}}$: First layer weights
- $\mathbf{b}_1 \in \mathbb{R}^h$: First layer biases
- σ : Activation function (nonlinearity)
- $\mathbf{h} \in \mathbb{R}^h$: Hidden representation
- $\mathbf{W}_2, \mathbf{b}_2$: Second layer parameters

[Universal Approximation Theorem] A single-hidden-layer MLP with enough neurons can approximate any continuous function on a compact domain to arbitrary accuracy.

Key insight: Width matters! But in practice, depth is often more efficient than width.

Intuition: Think of hidden neurons as "basis functions." With enough of them, you can represent complex functions as combinations of simple ones.

6.2.2 Depth vs Width

Wide networks (few layers, many neurons per layer):

- Pro: Easier to optimize (fewer layers = less gradient propagation)
- Pro: Universal approximation holds
- Con: Need exponentially many neurons for complex functions
- Con: Don't learn hierarchical features

Deep networks (many layers, moderate neurons per layer):

- Pro: More parameter efficient (exponentially fewer parameters)
- Pro: Learn hierarchical representations (low-level \rightarrow high-level features)
- Pro: Better generalization (implicit regularization)
- Con: Harder to optimize (vanishing/exploding gradients)
- Con: Requires careful initialization and normalization

Rule of thumb:

- Start with 2-3 hidden layers
- Width: 1-2 \times input dimension for first layer, then decrease
- Increase depth if data is complex and you have lots of samples
- Use modern techniques (batch norm, skip connections) for very deep networks

6.2.3 Activation Functions

Activation functions introduce nonlinearity. Without them, multiple layers would collapse to a single linear transformation.

ReLU (Rectified Linear Unit):

$$\text{ReLU}(x) = \max(0, x)$$

- **Pros:** Fast to compute, sparse activations, no vanishing gradient for $x > 0$
- **Cons:** "Dead ReLUs" (neurons that always output 0)
- **Derivative:** $\frac{d}{dx}\text{ReLU}(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$
- **Use when:** Default choice for most tasks, especially deep networks

Leaky ReLU:

$$\text{LeakyReLU}(x) = \max(\alpha x, x) \quad (\text{typically } \alpha = 0.01)$$

- **Pros:** Fixes dead ReLU problem (always has gradient)
- **Cons:** Extra hyperparameter α
- **Use when:** ReLU causes many dead neurons

ELU (Exponential Linear Unit):

$$\text{ELU}(x) = \begin{cases} x & x > 0 \\ \alpha(e^x - 1) & x \leq 0 \end{cases}$$

- **Pros:** Smooth, mean activations closer to zero
- **Cons:** Slower to compute (exponential)
- **Use when:** You want smoother gradients

Tanh (Hyperbolic Tangent):

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- **Pros:** Zero-centered (outputs in $[-1, 1]$), smooth
- **Cons:** Vanishing gradient for large $|x|$
- **Use when:** Need zero-centered activations (RNNs, some scientific applications)

Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- **Pros:** Outputs in $(0, 1)$ (interpretable as probabilities)
- **Cons:** Strong vanishing gradient, not zero-centered
- **Use when:** **Only in output layer** for binary classification

GELU (Gaussian Error Linear Unit):

$$\text{GELU}(x) = x \cdot \Phi(x) \quad \text{where } \Phi \text{ is CDF of } \mathcal{N}(0, 1)$$

- **Pros:** Smooth, used in state-of-the-art models (BERT, GPT)
- **Cons:** More expensive than ReLU
- **Use when:** Training Transformers or want best performance

Swish (SiLU):

$$\text{Swish}(x) = x \cdot \sigma(x)$$

- **Pros:** Smooth, unbounded above, self-gated
- **Cons:** More expensive than ReLU
- **Use when:** Similar to GELU, modern architecture

[Choosing Activation Functions] **Default choice:** Use **ReLU** for hidden layers unless you have a reason not to.

Use Leaky ReLU or ELU if: You observe many dead neurons (check activation distributions)

Use Tanh if: You need zero-centered activations (RNNs, specific scientific applications)

Use GELU or Swish if: You're training a Transformer or want state-of-the-art performance and can afford the compute

Never use Sigmoid in hidden layers: Strong vanishing gradient makes training very difficult

6.2.4 Why Networks Fail to Train

1. Vanishing Gradients

As gradients backpropagate through many layers, they can become exponentially small.

Cause: Chain rule multiplies many small derivatives (< 1)

Effect: Early layers learn very slowly or not at all

Solutions:

- Use ReLU instead of sigmoid/tanh (ReLU has gradient 1 for $x > 0$)
- Proper initialization (Xavier, He)
- Batch normalization
- Skip connections (ResNets—we'll cover this later)

2. Exploding Gradients

Gradients become exponentially large, causing NaN or very large parameter updates.

Cause: Chain rule multiplies many large derivatives (> 1)

Effect: Training diverges, loss becomes NaN

Solutions:

- Gradient clipping
- Proper initialization
- Lower learning rate
- Batch normalization

3. Dead ReLUs

ReLU neurons that always output 0 (because input is always negative).

Cause: Large negative bias or poor initialization

Effect: Neurons become permanently inactive, reducing model capacity

Detection:

```
1 # Check percentage of dead neurons
2 activations = model(x)
3 dead_percentage = (activations == 0).float().mean()
4 print(f"Dead neurons: {dead_percentage*100:.2f}%")
```

Solutions:

- Use Leaky ReLU or ELU
- Proper initialization (He initialization)
- Lower learning rate
- Batch normalization

4. Poor Initialization

Starting with inappropriate parameter values.

Effects:

- Too small: Vanishing activations/gradients
- Too large: Exploding activations/gradients

- All same: All neurons learn same features (symmetry problem)

Solution: Use proper initialization schemes (next section).

6.3 Implementation: Building MLPs in PyTorch

6.3.1 Simple MLP with Sequential

```

1 import torch
2 import torch.nn as nn
3
4 # 3-layer MLP: input_dim -> 128 -> 64 -> output_dim
5 model = nn.Sequential(
6     nn.Linear(10, 128),
7     nn.ReLU(),
8     nn.Linear(128, 64),
9     nn.ReLU(),
10    nn.Linear(64, 1)
11 )
12
13 # Test
14 x = torch.randn(32, 10) # Batch of 32 samples
15 y = model(x)
16 print(y.shape) # torch.Size([32, 1])

```

6.3.2 Custom MLP Module

```

1 class MLP(nn.Module):
2     """Flexible MLP with configurable architecture."""
3
4     def __init__(self, input_dim, hidden_dims, output_dim,
5                 activation='relu', dropout=0.0):
6         """
7         Args:
8             input_dim: Input feature dimension
9             hidden_dims: List of hidden layer dimensions
10            output_dim: Output dimension
11            activation: Activation function ('relu', 'tanh',
12            'gelu', etc.)
13            dropout: Dropout probability (0 = no dropout)
14        """
15        super().__init__()
16
17        # Build layers
18        layers = []
19        prev_dim = input_dim
20
21        for hidden_dim in hidden_dims:
22            layers.append(nn.Linear(prev_dim, hidden_dim))
23
24            # Activation
25            if activation == 'relu':
26                layers.append(nn.ReLU())
27            elif activation == 'tanh':
28                layers.append(nn.Tanh())
29            elif activation == 'gelu':
30                layers.append(nn.GELU())
31            elif activation == 'leaky_relu':
32                layers.append(nn.LeakyReLU(0.01))
33
34            # Dropout (if specified)

```

```
34         if dropout > 0:
35             layers.append(nn.Dropout(dropout))
36
37         prev_dim = hidden_dim
38
39         # Output layer (no activation)
40         layers.append(nn.Linear(prev_dim, output_dim))
41
42         self.network = nn.Sequential(*layers)
43
44     def forward(self, x):
45         return self.network(x)
46
47     # Usage
48     model = MLP(input_dim=10, hidden_dims=[128, 64, 32],
49                 output_dim=1,
50                 activation='relu', dropout=0.2)
51
52     print(model)
53     """
54     MLP(
55       (network): Sequential(
56         (0): Linear(in_features=10, out_features=128, bias=True)
57         (1): ReLU()
58         (2): Dropout(p=0.2, inplace=False)
59         (3): Linear(in_features=128, out_features=64, bias=True)
60         (4): ReLU()
61         (5): Dropout(p=0.2, inplace=False)
62         ...
63       )
64     )
65     """
```

6.3.3 Weight Initialization

Why initialization matters:

```

1 # Bad initialization (too small)
2 model = nn.Linear(100, 100)
3 with torch.no_grad():
4     model.weight.fill_(0.01)
5     model.bias.fill_(0)
6
7 x = torch.randn(1, 100)
8 for i in range(10):
9     x = torch.relu(model(x))
10    print(f"Layer {i}: mean={x.mean():.4f}, std={x.std():.4f}")
11 # Output: Values shrink to near zero (vanishing)
12
13 # Bad initialization (too large)
14 model = nn.Linear(100, 100)
15 with torch.no_grad():
16     model.weight.fill_(1.0)
17     model.bias.fill_(0)
18
19 x = torch.randn(1, 100)
20 for i in range(10):
21     x = torch.relu(model(x))
22     print(f"Layer {i}: mean={x.mean():.4f}, std={x.std():.4f}")
23 # Output: Values explode (exploding)

```

Xavier/Glorot Initialization (for Tanh/Sigmoid):

Maintains variance across layers for linear activations.

$$W \sim \mathcal{U}\left(-\sqrt{\frac{6}{n_{in} + n_{out}}}, \sqrt{\frac{6}{n_{in} + n_{out}}}\right)$$

```

1 import torch.nn.init as init
2
3 # Xavier uniform
4 init.xavier_uniform_(model.weight)
5
6 # Xavier normal
7 init.xavier_normal_(model.weight)

```

He/Kaiming Initialization (for ReLU):

Accounts for ReLU killing half the neurons (setting them to 0).

$$W \sim \mathcal{N}\left(0, \sqrt{\frac{2}{n_{in}}}\right)$$

```

1 # He/Kaiming normal (preferred for ReLU)
2 init.kaiming_normal_(model.weight, mode='fan_in',
3                       nonlinearity='relu')
4
5 # He/Kaiming uniform
6 init.kaiming_uniform_(model.weight, nonlinearity='relu')

```

Applying initialization to your model:

```

1 def init_weights(m):
2     """Initialize weights for Linear layers."""
3     if isinstance(m, nn.Linear):
4         # He initialization for weights
5         init.kaiming_normal_(m.weight, mode='fan_in',
6                               nonlinearity='relu')
7         # Zero initialization for biases
8         if m.bias is not None:
9             init.constant_(m.bias, 0)
10
11 # Apply to model
12 model = MLP(10, [128, 64], 1)
13 model.apply(init_weights)
14
15 # Verify
16 for name, param in model.named_parameters():
17     if 'weight' in name:
18         print(f"{name}: mean={param.mean():.4f}, std={param.
19               std():.4f}")

```

[Initialization Guidelines] **For ReLU networks:** Use He/Kaiming initialization

```
1 init.kaiming_normal_(weight, nonlinearity='relu')
```

For Tanh/Sigmoid networks: Use Xavier/Glorot initialization

```
1 init.xavier_normal_(weight)
```

For biases: Initialize to zero (or small constant)

```
1 init.constant_(bias, 0)
```

PyTorch defaults: Most layers use reasonable defaults (Kaiming uniform for Linear). Explicit initialization is often unnecessary but can help for deep networks.

6.3.4 Batch Normalization

Normalize activations within each mini-batch to stabilize training.

How it works:

$$\hat{x} = \frac{x - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y = \gamma \hat{x} + \beta$$

where μ_B , σ_B are batch mean/std, and γ , β are learnable parameters.

```

1 class MLPWithBatchNorm(nn.Module):
2     def __init__(self, input_dim, hidden_dims, output_dim):
3         super().__init__()
4
5         layers = []
6         prev_dim = input_dim
7
8         for hidden_dim in hidden_dims:
9             layers.append(nn.Linear(prev_dim, hidden_dim))
10            layers.append(nn.BatchNorm1d(hidden_dim)) # Add
11            # BatchNorm
12            layers.append(nn.ReLU())
13            prev_dim = hidden_dim
14
15            layers.append(nn.Linear(prev_dim, output_dim))
16            self.network = nn.Sequential(*layers)
17
18        def forward(self, x):
19            return self.network(x)
20
21    # Usage
22    model = MLPWithBatchNorm(10, [128, 64], 1)

```

Where to place BatchNorm?

Two conventions:

1. Linear \rightarrow BatchNorm \rightarrow Activation (more common)
2. Linear \rightarrow Activation \rightarrow BatchNorm

Both work; the first is more standard.

Training vs Eval mode:

```

1 model.train() # Use batch statistics
2 output = model(x)
3
4 model.eval() # Use running statistics
5 with torch.no_grad():
6     output = model(x)

```

[BatchNorm Pitfalls]

1. **Small batches:** BatchNorm works poorly with `batch_size < 8`. Use LayerNorm instead.
2. **Forgetting eval():** Model behaves differently in train/eval mode
3. **Before or after activation:** Be consistent in your architecture

6.3.5 Layer Normalization

Normalize across features instead of batch:

$$\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

where μ, σ are computed per sample across features.

```

1 class MLPWithLayerNorm(nn.Module):
2     def __init__(self, input_dim, hidden_dims, output_dim):
3         super().__init__()
4
5         layers = []
6         prev_dim = input_dim
7
8         for hidden_dim in hidden_dims:
9             layers.append(nn.Linear(prev_dim, hidden_dim))
10            layers.append(nn.LayerNorm(hidden_dim)) #
11            LayerNorm
12            layers.append(nn.ReLU())
13            prev_dim = hidden_dim
14
15            layers.append(nn.Linear(prev_dim, output_dim))
16            self.network = nn.Sequential(*layers)
17
18        def forward(self, x):
19            return self.network(x)

```

BatchNorm vs LayerNorm:

Aspect	BatchNorm	LayerNorm
Normalizes over	Batch dimension	Feature dimension
Requires batches	Yes (issues with batch=1)	No
Train/eval modes	Different	Same
Use in	CNNs, MLPs	Transformers, RNNs
Typical batch size	≥ 8	Any (even 1)

Table 3: BatchNorm vs LayerNorm comparison

6.3.6 Dropout

Randomly zero out neurons during training to prevent overfitting:

```

1 class MLPWithDropout(nn.Module):
2     def __init__(self, input_dim, hidden_dims, output_dim,
3         dropout=0.5):
4         super().__init__()
5         layers = []
6         prev_dim = input_dim
7
8         for hidden_dim in hidden_dims:
9             layers.append(nn.Linear(prev_dim, hidden_dim))
10            layers.append(nn.ReLU())
11            layers.append(nn.Dropout(dropout)) # Dropout
12            # after activation
13            prev_dim = hidden_dim
14
15            layers.append(nn.Linear(prev_dim, output_dim))
16            self.network = nn.Sequential(*layers)
17
18            def forward(self, x):
19                return self.network(x)
20
21            # Usage
22            model = MLPWithDropout(10, [128, 64], 1, dropout=0.5)
23
24            # Training: dropout active
25            model.train()
26            output = model(x) # Some neurons randomly zeroed
27
28            # Eval: dropout disabled
29            model.eval()
30            output = model(x) # All neurons active, outputs scaled

```

Dropout probabilities:

- Typical values: 0.2 - 0.5
- Higher dropout = more regularization (but too high hurts performance)
- Often use different dropout rates per layer (higher in later layers)

Where to place dropout:

- After activation functions
- **Not** before the output layer (usually)
- Can place after input layer (input dropout, typically lower rate like 0.1)

6.3.7 Complete MLP with All Techniques

```

1 class ModernMLP(nn.Module):
2     """MLP with all modern techniques."""
3
4     def __init__(self, input_dim, hidden_dims, output_dim,
5         activation='relu', dropout=0.0,
6         use_batch_norm=True):
7         super().__init__()
8
9         self.input_bn = nn.BatchNorm1d(input_dim) if
10            use_batch_norm else None

```

```

9
10     layers = []
11     prev_dim = input_dim
12
13     for i, hidden_dim in enumerate(hidden_dims):
14         # Linear layer
15         layers.append(nn.Linear(prev_dim, hidden_dim))
16
17         # Batch normalization
18         if use_batch_norm:
19             layers.append(nn.BatchNorm1d(hidden_dim))
20
21         # Activation
22         if activation == 'relu':
23             layers.append(nn.ReLU())
24         elif activation == 'gelu':
25             layers.append(nn.GELU())
26         elif activation == 'leaky_relu':
27             layers.append(nn.LeakyReLU(0.01))
28
29         # Dropout (if specified)
30         if dropout > 0:
31             layers.append(nn.Dropout(dropout))
32
33         prev_dim = hidden_dim
34
35     # Output layer
36     layers.append(nn.Linear(prev_dim, output_dim))
37
38     self.network = nn.Sequential(*layers)
39
40     # Initialize weights
41     self.apply(self._init_weights)
42
43     def _init_weights(self, m):
44         if isinstance(m, nn.Linear):
45             init.kaiming_normal_(m.weight, nonlinearity='relu
46         ')
47             if m.bias is not None:
48                 init.constant_(m.bias, 0)
49
50     def forward(self, x):
51         # Optional input batch norm
52         if self.input_bn is not None:
53             x = self.input_bn(x)
54         return self.network(x)
55
56 # Usage: best practices
57 model = ModernMLP(
58     input_dim=10,
59     hidden_dims=[256, 128, 64],
60     output_dim=1,
61     activation='relu',
62     dropout=0.3,
63     use_batch_norm=True
64 )

```

6.4 Implementation: Debugging MLPs

6.4.1 Checking for Dead ReLUs

```

1 def check_dead_neurons(model, dataloader, device):
2     """Check percentage of dead ReLU neurons."""
3     model.eval()
4
5     # Hooks to capture activations
6     activations = {}
7
8     def hook_fn(name):
9         def hook(module, input, output):
10             if name not in activations:
11                 activations[name] = []
12                 activations[name].append((output == 0).float())
13             return hook
14
15     # Register hooks on ReLU layers
16     hooks = []
17     for name, module in model.named_modules():
18         if isinstance(module, nn.ReLU):
19             hooks.append(module.register_forward_hook(hook_fn
20 (name)))
21
22     # Forward pass on some data
23     with torch.no_grad():
24         for data, _ in dataloader:
25             data = data.to(device)
26             model(data)
27             break # Just one batch
28
29     # Remove hooks
30     for hook in hooks:
31         hook.remove()
32
33     # Compute dead neuron percentage
34     for name, acts in activations.items():
35         dead_pct = torch.cat(acts).mean().item() * 100
36         print(f"{name}: {dead_pct:.2f}% neurons always zero")
37
38 # Usage
39 # check_dead_neurons(model, train_loader, device)

```

6.4.2 Visualizing Activation Distributions

```

1 import matplotlib.pyplot as plt
2
3 def plot_activation_distributions(model, x):
4     """Plot distribution of activations in each layer."""
5     model.eval()
6
7     activations = []
8
9     def hook_fn(module, input, output):
10         activations.append(output.detach().cpu().flatten())
11
12     # Register hooks
13     hooks = []
14     for name, module in model.named_modules():
15         if isinstance(module, (nn.Linear, nn.ReLU)):

```

```

16         hooks.append(module.register_forward_hook(hook_fn
17     ))
18     # Forward pass
19     with torch.no_grad():
20         model(x)
21
22     # Remove hooks
23     for hook in hooks:
24         hook.remove()
25
26     # Plot
27     fig, axes = plt.subplots(len(activations), 1,
28                             figsize=(10, 3*len(activations)))
29     for i, act in enumerate(activations):
30         axes[i].hist(act.numpy(), bins=50)
31         axes[i].set_title(f'Layer {i} activations')
32         axes[i].set_xlabel('Value')
33         axes[i].set_ylabel('Count')
34     plt.tight_layout()
35     plt.show()
36
37 # Usage
38 # x = torch.randn(100, 10)
39 # plot_activation_distributions(model, x)

```

6.4.3 Gradient Flow Visualization

```

1 def plot_gradient_flow(model):
2     """
3     Plot gradient flow through network.
4     Call after loss.backward() but before optimizer.step()
5     """
6     ave_grads = []
7     max_grads = []
8     layers = []
9
10    for name, param in model.named_parameters():
11        if param.requires_grad and param.grad is not None:
12            layers.append(name)
13            ave_grads.append(param.grad.abs().mean().item())
14            max_grads.append(param.grad.abs().max().item())
15
16    plt.figure(figsize=(12, 6))
17    plt.bar(range(len(ave_grads)), ave_grads, alpha=0.5,
18            label='mean')
19    plt.bar(range(len(max_grads)), max_grads, alpha=0.5,
20            label='max')
21    plt.hlines(0, 0, len(ave_grads), linewidth=2, color='k')
22    plt.xticks(range(len(layers)), layers, rotation='vertical',
23              ')
24    plt.xlim(-1, len(ave_grads))
25    plt.xlabel('Layers')
26    plt.ylabel('Gradient magnitude')
27    plt.legend()
28    plt.title('Gradient Flow')
29    plt.grid(True)
30    plt.tight_layout()
31    plt.show()
32
33 # Usage in training loop

```

```
31 # output = model(x)
32 # loss = criterion(output, y)
33 # loss.backward()
34 # plot_gradient_flow(model) # Before optimizer.step()
35 # optimizer.step()
```

6.5 Exercises

6.1: Building Your First MLP - ★★

Goal: Create and train a basic MLP.

1. Build an MLP: $20 \rightarrow 64 \rightarrow 32 \rightarrow 1$
2. Use ReLU activations
3. Train on synthetic data: $y = \sum_{i=1}^{20} x_i + \epsilon$
4. Use 1000 samples, MSE loss, Adam optimizer
5. Train for 100 epochs
6. Plot training loss curve

Success criterion: Final loss < 0.1

Starter code:

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4
5 # Generate data
6 X = torch.randn(1000, 20)
7 y = X.sum(dim=1, keepdim=True) + 0.1 * torch.randn(1000,
8     1)
9
10 # Your MLP here
11 model = ...
12 # Training loop here

```

6.2: Activation Function Comparison - ★★★

Goal: Understand how different activations affect training.

Train identical MLPs (3 layers, 64 hidden units) with different activations:

1. ReLU
2. Leaky ReLU
3. Tanh
4. GELU
5. Sigmoid (observe why this is bad for hidden layers!)

For each:

- Train on the same task (e.g., nonlinear function approximation)
- Track training loss over epochs
- Measure final test performance
- Plot all loss curves on the same graph

Questions to answer:

- Which converges fastest?
- Which achieves best final performance?
- Why does sigmoid perform poorly?

6.3: Initialization Impact - ★★★

Goal: See how initialization affects training.

Create the same MLP with different initializations:

1. All zeros (model won't learn—why?)
2. All ones (all neurons learn the same features)
3. Random normal with $\sigma = 0.01$ (too small)
4. Random normal with $\sigma = 1.0$ (too large)
5. Xavier initialization
6. He initialization (best for ReLU)

For each:

- Train for 50 epochs
- Track loss
- Check for NaN values

Visualize: Plot all loss curves. Which initializations fail? Which work best?

6.4: Depth vs Width Experiment - ★★

Goal: Compare deep vs wide networks.

Create networks with approximately the same number of parameters:

1. **Shallow & wide:** 2 layers, 256 neurons each
2. **Deep & narrow:** 8 layers, 64 neurons each
3. **Medium:** 4 layers, 128 neurons each

Train all three on a moderately complex task (e.g., 2D function approximation).

Compare:

- Training speed (epochs to converge)
- Final test performance
- Generalization (train vs test loss gap)

Challenge: Add batch normalization to the deep network. Does it help?

6.5: Regularization Techniques - ★★

Goal: Prevent overfitting with regularization.

1. Create a small dataset (200 samples)
2. Build a large MLP (can easily overfit): $10 \rightarrow 128 \rightarrow 128 \rightarrow 64 \rightarrow 1$
3. Train without regularization (observe overfitting)
4. Add dropout (try 0.1, 0.3, 0.5)
5. Add batch normalization
6. Add weight decay
7. Combine techniques

For each variant:

- Track train and validation loss
- Measure final test performance
- Plot train vs val loss curves

Questions:

- Which technique helps most?
- What's the best combination?
- Can you achieve zero overfitting gap?

6.6: Complex Function Approximation - ★★★★★**Goal:** Use MLP for scientific computing task.

Approximate a complex 2D function:

$$f(x, y) = \sin(5x) \cos(5y) + 0.5 \sin(10x) + 0.3 \cos(10y)$$

1. Generate training data: sample (x, y) uniformly in $[-1, 1]^2$
2. Create MLP: 2 inputs \rightarrow hidden layers \rightarrow 1 output
3. Train with MSE loss
4. Visualize:
 - True function as heatmap
 - MLP prediction as heatmap
 - Absolute error heatmap
5. Try different architectures (depth, width)
6. Try different activations

Success criterion: Mean absolute error < 0.05 on test set**Starter code for visualization:**

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Create grid
5 x = np.linspace(-1, 1, 100)
6 y = np.linspace(-1, 1, 100)
7 X, Y = np.meshgrid(x, y)
8
9 # True function
10 Z_true = np.sin(5*X)*np.cos(5*Y) + 0.5*np.sin(10*X) +
    0.3*np.cos(10*Y)
11
12 # Model prediction
13 grid_points = torch.FloatTensor(np.stack([X.ravel(), Y.
    ravel()], axis=1))
14 with torch.no_grad():
15     Z_pred = model(grid_points).numpy().reshape(X.shape)
16
17 # Plot
18 fig, axes = plt.subplots(1, 3, figsize=(15, 4))
19 axes[0].contourf(X, Y, Z_true, levels=20)
20 axes[0].set_title('True Function')
21 axes[1].contourf(X, Y, Z_pred, levels=20)
22 axes[1].set_title('MLP Prediction')
23 axes[2].contourf(X, Y, np.abs(Z_true - Z_pred), levels
    =20)
24 axes[2].set_title('Absolute Error')
25 plt.show()

```

6.6 Key Takeaways

MLPs are versatile:

- Universal function approximators
- Foundation for understanding all neural networks
- Effective for tabular data and function approximation

Design choices matter:

- **Activation:** ReLU is default, but consider alternatives for specific tasks
- **Depth vs Width:** Deeper networks more parameter-efficient but harder to train
- **Initialization:** He for ReLU, Xavier for Tanh, crucial for deep networks
- **Normalization:** BatchNorm stabilizes training, especially for deep networks
- **Regularization:** Dropout and weight decay prevent overfitting

Common failure modes:

- Dead ReLUs → Use Leaky ReLU or proper initialization
- Vanishing gradients → Use ReLU, batch norm, or skip connections
- Exploding gradients → Use gradient clipping, lower LR, batch norm
- Overfitting → Add dropout, weight decay, or get more data

Best practices for MLPs:

1. Start with 2-3 hidden layers
2. Use ReLU activation
3. Initialize with He initialization
4. Add batch normalization for deep networks
5. Use dropout (0.2-0.5) to prevent overfitting
6. Monitor both train and validation loss
7. Visualize activations and gradients when debugging

7 Convolutional Neural Networks (CNNs)

7.1 Introduction: Why Convolutions?

MLPs work well for tabular data, but they have serious limitations for spatial data like images:

Problems with MLPs for images:

- **Too many parameters:** A 224×224 RGB image has 150,528 inputs. First hidden layer with 1000 neurons = 150M parameters!
- **No spatial structure:** Treats each pixel independently, ignoring that nearby pixels are related
- **Not translation invariant:** A cat in the top-left corner is different from a cat in the center
- **Can't generalize:** Learning a feature at one position doesn't help recognize it elsewhere

Convolutions solve these problems through:

1. **Parameter sharing:** Same kernel applied everywhere (drastically fewer parameters)
2. **Local connectivity:** Each neuron only looks at a small region (receptive field)
3. **Translation invariance:** Features learned at one position work everywhere
4. **Hierarchical learning:** Early layers detect edges, later layers detect complex patterns

When to use CNNs:

- Images (most common use case)
- Time series (1D convolutions)
- Volumetric data (3D medical scans, physical simulations)
- Spatial/temporal data with local correlations

When NOT to use CNNs:

- Tabular data (no spatial structure)
- Very small spatial dimensions (convolutions overkill)
- When you need long-range dependencies only (use attention/Transformers)

7.2 Theory: How Convolutions Work

7.2.1 The Convolution Operation

A 2D convolution slides a **kernel** (or **filter**) over the input, computing dot products:

$$\text{Output}[i, j] = \sum_{m=0}^{k-1} \sum_{n=0}^{k-1} \text{Input}[i + m, j + n] \cdot \text{Kernel}[m, n]$$

Example: 3×3 kernel on 5×5 input

Input (5x5):		Kernel (3x3):
1 2 3 4 5		1 0 -1
2 3 4 5 6	*	2 0 -2
3 4 5 6 7		1 0 -1

```
4 5 6 7 8
5 6 7 8 9
```

Output (3x3):

```
-4 -4 -4
-4 -4 -4
-4 -4 -4
```

This particular kernel is a vertical edge detector (Sobel filter).

[Key Concepts]

- **Kernel/Filter:** Small matrix of learnable weights (e.g., 3×3 , 5×5)
- **Stride:** How many pixels to move the kernel each step (stride=1 \rightarrow move 1 pixel)
- **Padding:** Add zeros around input to control output size
- **Channels:** Input depth (RGB image has 3 channels)
- **Feature maps:** Output of applying convolution (number = number of kernels)

7.2.2 Parameter Sharing

A single 3×3 kernel has only 9 parameters, but is applied to every position in the image.

MLP for 28×28 image:

- Input: 784 pixels
- Hidden layer: 128 neurons
- Parameters: $784 \times 128 = 100,352$

CNN for same image:

- Input: $28 \times 28 \times 1$
- Conv layer: 32 filters, 3×3 kernel
- Parameters: $3 \times 3 \times 1 \times 32 + 32 = 320$ (including biases)

300 \times fewer parameters!

7.2.3 Receptive Fields

The **receptive field** of a neuron is the region of the input it "sees."

Example: Two 3×3 conv layers

Layer 1: Each neuron sees 3×3 region

Layer 2: Each neuron sees 3×3 of Layer 1

→ Sees 5×5 region of original input!

Receptive field grows with depth:

- 1 layer (3×3 kernel): $3 \times 3 = 9$ pixels
- 2 layers: $5 \times 5 = 25$ pixels
- 3 layers: $7 \times 7 = 49$ pixels
- n layers: $(2n + 1) \times (2n + 1)$ pixels (for 3×3 kernels, stride=1)

Formula for receptive field:

$$r_n = r_{n-1} + (k - 1) \cdot \prod_{i=1}^{n-1} s_i$$

where k is kernel size, s_i is stride at layer i .

[Why Receptive Fields Matter] For a network to classify an entire image, the final layer's receptive field must cover the whole image. This determines your network depth:

Small images (32×32): 5-10 layers sufficient

Large images (224×224): 15-20+ layers needed

Use pooling or larger strides to grow receptive field faster.

7.2.4 Output Size Calculation

Formula for output size:

$$\text{Output size} = \left\lfloor \frac{\text{Input size} + 2 \times \text{Padding} - \text{Kernel size}}{\text{Stride}} \right\rfloor + 1$$

Examples:

Input	Kernel	Stride	Padding	Output	Note
28	3	1	0	26	Shrinks by 2
28	3	1	1	28	Same size
28	3	2	0	13	Half size
28	5	1	2	28	Same size
32	3	1	0	30	Shrinks by 2

Table 4: Output size examples

Common padding strategies:

- **Valid (no padding):** Output shrinks
- **Same:** Output size = input size (for stride=1)
- **Full:** Maximum padding (output grows)

7.2.5 Pooling

Pooling **downsamples** the spatial dimensions, making the network:

- More computationally efficient
- More translation invariant
- With larger receptive fields

Max Pooling: Take maximum value in each window

Input (4x4):		Max Pool 2x2, stride=2:
1 2 3 4		
2 4 6 8	→	4 8
3 6 9 12		12 16
4 8 12 16		

Average Pooling: Take average value in each window

Input (4x4):		Avg Pool 2x2, stride=2:
1 2 3 4		
2 4 6 8	→	2.25 5.25
3 6 9 12		5.25 11.25
4 8 12 16		

When to use each:

- **Max pooling:** Most common. Preserves strong features. Good for classification.
- **Average pooling:** Smoother. Good for final global pooling. Sometimes better for regression.
- **No pooling:** Use strided convolutions instead (more learnable, modern approach)

7.3 Implementation: Building CNNs in PyTorch

7.3.1 Conv2d: The Core Layer

```

1 import torch
2 import torch.nn as nn
3
4 # Basic Conv2d
5 conv = nn.Conv2d(
6     in_channels=3,          # Input channels (e.g., RGB = 3)
7     out_channels=64,        # Number of filters (output channels)
8     kernel_size=3,         # 3x3 kernel
9     stride=1,              # Move 1 pixel at a time
10    padding=1,              # Pad with 1 pixel of zeros (same
11                             size)
12    bias=True               # Include bias term
13 )
14 # Input: (batch_size, channels, height, width)
15 x = torch.randn(8, 3, 32, 32) # 8 RGB images, 32x32
16 output = conv(x)
17 print(output.shape) # torch.Size([8, 64, 32, 32])
18 # 64 feature maps, same spatial size due to padding=1
19
20 # Number of parameters
21 params = conv.weight.numel() + conv.bias.numel()
22 print(f"Parameters: {params}") # 3*3*3*64 + 64 = 1792

```

Key parameters explained:

- **in_channels:** Depth of input (3 for RGB, 1 for grayscale)
- **out_channels:** Number of filters = depth of output
- **kernel_size:** Can be int ($3 \rightarrow 3 \times 3$) or tuple ($3, 5 \rightarrow 3 \times 5$)
- **stride:** Step size. stride=2 halves spatial dimensions
- **padding:**
 - int: uniform padding
 - tuple: (pad_height, pad_width)
 - 'same': auto-pad to keep size (stride=1 only)
 - 'valid': no padding
- **dilation:** Spacing between kernel elements (1 = standard, >1 = dilated)
- **groups:** Divide channels into groups (1 = standard, in_channels = depthwise)

7.3.2 Shape Calculation Practice

```

1 # Practice calculating output shapes
2 x = torch.randn(1, 3, 64, 64) # Single 64x64 RGB image
3
4 # Example 1: Standard convolution
5 conv1 = nn.Conv2d(3, 32, kernel_size=3, padding=1)
6 out1 = conv1(x)
7 print(out1.shape) # [1, 32, 64, 64] - same size
8
9 # Example 2: Strided convolution (downsampling)
10 conv2 = nn.Conv2d(3, 32, kernel_size=3, stride=2, padding=1)
11 out2 = conv2(x)
12 print(out2.shape) # [1, 32, 32, 32] - halved
13
14 # Example 3: No padding (shrinks)

```

```

15 conv3 = nn.Conv2d(3, 32, kernel_size=5, padding=0)
16 out3 = conv3(x)
17 print(out3.shape)  # [1, 32, 60, 60] - shrinks by 4
18
19 # Example 4: Large kernel
20 conv4 = nn.Conv2d(3, 32, kernel_size=7, padding=3)
21 out4 = conv4(x)
22 print(out4.shape)  # [1, 32, 64, 64] - same size
23
24 # Calculate expected output
25 def calc_output_size(input_size, kernel_size, stride, padding
26 ):
27     return (input_size + 2*padding - kernel_size) // stride +
28         1
29
30 print(calc_output_size(64, 3, 1, 1))  # 64
31 print(calc_output_size(64, 3, 2, 1))  # 32
32 print(calc_output_size(64, 5, 1, 0))  # 60

```

[Common Shape Mistakes] **Mistake 1:** Forgetting to flatten before Linear layer

```

1 # WRONG
2 output = conv_layers(x)  # Shape: (batch, channels, H, W)
3 output = fc(output)      # ERROR! Wrong shape
4
5 # CORRECT
6 output = conv_layers(x)
7 output = output.view(output.size(0), -1)  # Flatten
8 output = fc(output)  # OK

```

Mistake 2: Wrong channel order (PyTorch uses NCHW, not NHWC)

```

1 # PyTorch expects: (batch, channels, height, width)
2 x = torch.randn(8, 32, 32, 3)  # WRONG order
3 x = x.permute(0, 3, 1, 2)      # Fix to (8, 3, 32, 32)

```


7.3.3 Building a Simple CNN

```

1 class SimpleCNN(nn.Module):
2     """Simple CNN for MNIST (28x28 grayscale images)."""
3
4     def __init__(self, num_classes=10):
5         super().__init__()
6
7         # Convolutional layers
8         self.conv1 = nn.Conv2d(1, 32, kernel_size=3, padding
9 =1)
10        # Output: 32 x 28 x 28
11
12        self.conv2 = nn.Conv2d(32, 64, kernel_size=3, padding
13 =1)
14        # Output: 64 x 28 x 28
15
16        self.pool = nn.MaxPool2d(2, 2) # 2x2 pooling
17        # After pool: 64 x 14 x 14
18
19        self.conv3 = nn.Conv2d(64, 128, kernel_size=3,
20 padding=1)
21        # Output: 128 x 14 x 14
22        # After pool: 128 x 7 x 7
23
24        # Fully connected layers
25        self.fc1 = nn.Linear(128 * 7 * 7, 256)
26        self.fc2 = nn.Linear(256, num_classes)
27
28        self.relu = nn.ReLU()
29        self.dropout = nn.Dropout(0.5)
30
31    def forward(self, x):
32        # Conv block 1
33        x = self.relu(self.conv1(x))
34        x = self.pool(x)
35
36        # Conv block 2
37        x = self.relu(self.conv2(x))
38        x = self.pool(x)
39
40        # Conv block 3
41        x = self.relu(self.conv3(x))
42
43        # Flatten
44        x = x.view(x.size(0), -1) # or x.flatten(1)
45
46        # Fully connected
47        x = self.relu(self.fc1(x))
48        x = self.dropout(x)
49        x = self.fc2(x)
50
51        return x
52
53    # Test
54    model = SimpleCNN()
55    x = torch.randn(4, 1, 28, 28) # 4 grayscale 28x28 images
56    output = model(x)
57    print(output.shape) # torch.Size([4, 10])

```

7.3.4 Modern CNN with Batch Normalization

```

1 class ModernCNN(nn.Module):
2     """CNN with batch normalization and modern practices."""
3
4     def __init__(self, num_classes=10):
5         super().__init__()
6
7         self.features = nn.Sequential(
8             # Block 1: 1 -> 32
9             nn.Conv2d(1, 32, 3, padding=1),
10            nn.BatchNorm2d(32),
11            nn.ReLU(),
12            nn.Conv2d(32, 32, 3, padding=1),
13            nn.BatchNorm2d(32),
14            nn.ReLU(),
15            nn.MaxPool2d(2, 2), # 28x28 -> 14x14
16
17            # Block 2: 32 -> 64
18            nn.Conv2d(32, 64, 3, padding=1),
19            nn.BatchNorm2d(64),
20            nn.ReLU(),
21            nn.Conv2d(64, 64, 3, padding=1),
22            nn.BatchNorm2d(64),
23            nn.ReLU(),
24            nn.MaxPool2d(2, 2), # 14x14 -> 7x7
25
26            # Block 3: 64 -> 128
27            nn.Conv2d(64, 128, 3, padding=1),
28            nn.BatchNorm2d(128),
29            nn.ReLU(),
30            nn.Conv2d(128, 128, 3, padding=1),
31            nn.BatchNorm2d(128),
32            nn.ReLU(),
33        )
34
35        self.classifier = nn.Sequential(
36            nn.AdaptiveAvgPool2d((1, 1)), # Global average
37            pooling
38            nn.Flatten(),
39            nn.Linear(128, 256),
40            nn.ReLU(),
41            nn.Dropout(0.5),
42            nn.Linear(256, num_classes)
43        )
44
45        def forward(self, x):
46            x = self.features(x)
47            x = self.classifier(x)
48            return x
49
50 # Test
51 model = ModernCNN()
52 x = torch.randn(4, 1, 28, 28)
53 output = model(x)
54 print(output.shape) # torch.Size([4, 10])

```

[AdaptiveAvgPool2d vs Flatten] **Old approach:** Fixed spatial size before FC layer

```

1 x = x.view(batch_size, -1) # Must know spatial size
2 fc = nn.Linear(128 * 7 * 7, num_classes)

```

Modern approach: Global pooling adapts to any input size

```
1 x = nn.AdaptiveAvgPool2d((1, 1))(x)  # Always outputs 1x1
2 x = x.flatten(1)  # (batch, 128, 1, 1) -> (batch, 128)
3 fc = nn.Linear(128, num_classes)
```

Benefit: Same network works for different input sizes (e.g., 28×28 or 32×32).

7.3.5 1D Convolutions for Time Series

```

1 class TimeSeriesCNN(nn.Module):
2     """1D CNN for time series or sequence data."""
3
4     def __init__(self, input_channels=1, num_classes=10):
5         super().__init__()
6
7         self.conv1 = nn.Conv1d(
8             in_channels=input_channels,
9             out_channels=64,
10            kernel_size=7, # Look at 7 time steps
11            padding=3
12        )
13
14        self.conv2 = nn.Conv1d(64, 128, kernel_size=5,
15                                padding=2)
16        self.conv3 = nn.Conv1d(128, 256, kernel_size=3,
17                                padding=1)
18
19        self.pool = nn.MaxPool1d(2) # Downsample by 2
20
21        self.global_pool = nn.AdaptiveAvgPool1d(1)
22        self.fc = nn.Linear(256, num_classes)
23
24        self.relu = nn.ReLU()
25
26    def forward(self, x):
27        # x shape: (batch, channels, sequence_length)
28
29        x = self.relu(self.conv1(x))
30        x = self.pool(x)
31
32        x = self.relu(self.conv2(x))
33        x = self.pool(x)
34
35        x = self.relu(self.conv3(x))
36
37        # Global pooling
38        x = self.global_pool(x) # (batch, 256, 1)
39        x = x.squeeze(-1) # (batch, 256)
40
41        x = self.fc(x)
42        return x
43
44    # Test on time series
45    model = TimeSeriesCNN(input_channels=1)
46    x = torch.randn(32, 1, 100) # 32 samples, 1 channel, 100
47                                     time steps
48    output = model(x)
49    print(output.shape) # torch.Size([32, 10])

```

When to use 1D convolutions:

- Time series data (sensor readings, stock prices)
- Audio signals (raw waveforms)
- Text (character or word level)
- Any 1D sequence with local correlations

7.3.6 3D Convolutions for Volumetric Data

```

1 class VolumetricCNN(nn.Module):
2     """3D CNN for volumetric data (medical scans, video, 3D
3     simulations)."""
4
5     def __init__(self, in_channels=1, num_classes=2):
6         super().__init__()
7
8         self.conv1 = nn.Conv3d(in_channels, 32, kernel_size
9         =3, padding=1)
10        self.conv2 = nn.Conv3d(32, 64, kernel_size=3, padding
11        =1)
12        self.conv3 = nn.Conv3d(64, 128, kernel_size=3,
13        padding=1)
14
15        self.pool = nn.MaxPool3d(2) # Downsample all 3
16        dimensions
17
18        self.global_pool = nn.AdaptiveAvgPool3d((1, 1, 1))
19        self.fc = nn.Linear(128, num_classes)
20
21        self.relu = nn.ReLU()
22
23    def forward(self, x):
24        # x shape: (batch, channels, depth, height, width)
25
26        x = self.relu(self.conv1(x))
27        x = self.pool(x)
28
29        x = self.relu(self.conv2(x))
30        x = self.pool(x)
31
32        x = self.relu(self.conv3(x))
33
34        x = self.global_pool(x)
35        x = x.view(x.size(0), -1)
36
37        x = self.fc(x)
38        return x
39
40 # Test on 3D volume
41 model = VolumetricCNN()
42 x = torch.randn(4, 1, 32, 32, 32) # 4 volumes, 32x32x32
43 output = model(x)
44 print(output.shape) # torch.Size([4, 2])

```

When to use 3D convolutions:

- Medical imaging (CT scans, MRI)
- Video (treat time as 3rd spatial dimension)
- 3D physical simulations
- Molecular structure analysis

Warning: 3D convolutions are very memory-intensive! Use smaller batch sizes.

7.3.7 Transposed Convolutions (Upsampling)

Transposed convolutions (sometimes called "deconvolutions") **increase** spatial dimensions. Used in:

- Autoencoders (decoder part)
- Generative models (GANs, VAEs)
- Semantic segmentation (U-Net)
- Super-resolution

```

1 # Regular convolution (downsampling)
2 conv = nn.Conv2d(64, 128, kernel_size=3, stride=2, padding=1)
3 x = torch.randn(1, 64, 32, 32)
4 out = conv(x)
5 print(out.shape) # torch.Size([1, 128, 16, 16]) - halved
6
7 # Transposed convolution (upsampling)
8 deconv = nn.ConvTranspose2d(128, 64, kernel_size=3, stride=2,
9                             padding=1,
10                             output_padding=1)
11 out2 = deconv(out)
12 print(out2.shape) # torch.Size([1, 64, 32, 32]) - doubled
13 back

```

Output size calculation:

Output size = (Input size - 1) × Stride - 2 × Padding + Kernel size + Output padding

Simple U-Net style autoencoder:

```

1 class SimpleAutoencoder(nn.Module):
2     """Convolutional autoencoder for images."""
3
4     def __init__(self):
5         super().__init__()
6
7         # Encoder (downsampling)
8         self.encoder = nn.Sequential(
9             nn.Conv2d(1, 32, 3, stride=2, padding=1), # 28
10            -> 14
11            nn.ReLU(),
12            nn.Conv2d(32, 64, 3, stride=2, padding=1), # 14
13            -> 7
14            nn.ReLU(),
15            nn.Conv2d(64, 128, 3, stride=2, padding=1), # 7
16            -> 4 (rounds up)
17            nn.ReLU(),
18        )
19
20        # Decoder (upsampling)
21        self.decoder = nn.Sequential(
22            nn.ConvTranspose2d(128, 64, 3, stride=2, padding
23            =1,
24            output_padding=1), # 4 -> 7
25            nn.ReLU(),
26            nn.ConvTranspose2d(64, 32, 3, stride=2, padding
27            =1,
28            output_padding=1), # 7 -> 14
29            nn.ReLU(),
30            nn.ConvTranspose2d(32, 1, 3, stride=2, padding=1,

```

```

26         output_padding=1), # 14 -> 28
27         nn.Sigmoid() # Output in [0, 1]
28     )
29
30     def forward(self, x):
31         encoded = self.encoder(x)
32         decoded = self.decoder(encoded)
33         return decoded
34
35     # Test
36     model = SimpleAutoencoder()
37     x = torch.randn(4, 1, 28, 28)
38     reconstructed = model(x)
39     print(reconstructed.shape) # torch.Size([4, 1, 28, 28])

```

[Checkerboard Artifacts] Transposed convolutions can create checkerboard artifacts (visible patterns in output).

Solution: Use `resize + convolution` instead:

```

1 # Instead of ConvTranspose2d
2 nn.ConvTranspose2d(64, 32, kernel_size=4, stride=2, padding
   =1)
3
4 # Use this (better quality)
5 nn.Sequential(
6     nn.Upsample(scale_factor=2, mode='bilinear',
7                 align_corners=False),
8     nn.Conv2d(64, 32, kernel_size=3, padding=1)
9 )

```

7.3.8 Dilated Convolutions

Dilated (atrous) convolutions increase receptive field without increasing parameters or losing resolution.

```

1 # Standard 3x3 convolution
2 conv_standard = nn.Conv2d(64, 64, kernel_size=3, padding=1)
3 # Receptive field: 3x3
4
5 # Dilated 3x3 convolution (dilation=2)
6 conv_dilated = nn.Conv2d(64, 64, kernel_size=3, padding=2,
   dilation=2)
7 # Receptive field: 5x5 (with gaps)
8 # Same parameters as standard 3x3!
9
10 # Effective kernel size: k + (k-1)*(d-1)
11 # For k=3, d=2: 3 + (3-1)*(2-1) = 5

```

When to use dilated convolutions:

- Semantic segmentation (maintain resolution while growing receptive field)
- Dense prediction tasks
- When you need large receptive field but want to keep resolution

7.3.9 Common CNN Architectures

VGG-style (stacking blocks):

```

1 def vgg_block(in_channels, out_channels, num_convs):
2     """VGG-style block: multiple convs + pool."""
3     layers = []
4     for _ in range(num_convs):
5         layers.append(nn.Conv2d(in_channels, out_channels,
6                                 kernel_size=3, padding=1))
7         layers.append(nn.ReLU())
8         in_channels = out_channels
9     layers.append(nn.MaxPool2d(2, 2))
10    return nn.Sequential(*layers)
11
12 class VGGStyleNet(nn.Module):
13     def __init__(self, num_classes=10):
14         super().__init__()
15
16         self.features = nn.Sequential(
17             vgg_block(3, 64, 2),      # 2 convs, 64 filters
18             vgg_block(64, 128, 2),    # 2 convs, 128 filters
19             vgg_block(128, 256, 3),    # 3 convs, 256 filters
20             vgg_block(256, 512, 3),    # 3 convs, 512 filters
21         )
22
23         self.classifier = nn.Sequential(
24             nn.AdaptiveAvgPool2d((1, 1)),
25             nn.Flatten(),
26             nn.Linear(512, 512),
27             nn.ReLU(),
28             nn.Dropout(0.5),
29             nn.Linear(512, num_classes)
30         )
31
32     def forward(self, x):
33         x = self.features(x)
34         x = self.classifier(x)
35         return x

```

Key design patterns:

1. **Double channels after pooling:** $64 \rightarrow 128 \rightarrow 256 \rightarrow 512$
2. **Multiple convs per block:** Learn richer features before downsampling
3. **Small kernels (3×3):** More efficient than large kernels
4. **Batch norm after conv:** Stabilizes training
5. **Global pooling:** More flexible than fixed FC input size

7.4 Implementation: Debugging CNNs

7.4.1 Visualizing Feature Maps

```

1 def visualize_feature_maps(model, x, layer_num=0):
2     """Visualize activations from a specific layer."""
3     activation = {}
4
5     def get_activation(name):
6         def hook(model, input, output):
7             activation[name] = output.detach()
8         return hook
9

```



```

10     # Register hook
11     layer_name = f'features.{layer_num}'
12     handle = dict(model.named_modules())[layer_name].
register_forward_hook(
13         get_activation(layer_name))
14
15     # Forward pass
16     model.eval()
17     with torch.no_grad():
18         _ = model(x)
19
20     handle.remove()
21
22     # Visualize
23     act = activation[layer_name].squeeze() # Remove batch
dim
24
25     import matplotlib.pyplot as plt
26     fig, axes = plt.subplots(4, 8, figsize=(16, 8))
27     for i, ax in enumerate(axes.flat):
28         if i < act.shape[0]:
29             ax.imshow(act[i].cpu(), cmap='viridis')
30             ax.axis('off')
31     plt.tight_layout()
32     plt.show()
33
34     # Usage
35     # model = ModernCNN()
36     # x = torch.randn(1, 1, 28, 28)
37     # visualize_feature_maps(model, x, layer_num=0)

```

7.4.2 Common CNN Debugging Issues

Problem: Output shape wrong

```

1  # Debug shape through network
2  def debug_shapes(model, input_shape):
3      """Print shapes through the network."""
4      x = torch.randn(*input_shape)
5      print(f"Input: {x.shape}")
6
7      for name, module in model.named_children():
8          x = module(x)
9          print(f"{name}: {x.shape}")
10
11     # Usage
12     model = SimpleCNN()
13     debug_shapes(model, (1, 1, 28, 28))
14     """
15     Input: torch.Size([1, 1, 28, 28])
16     conv1: torch.Size([1, 32, 28, 28])
17     pool: torch.Size([1, 32, 14, 14])
18     ...
19     """

```

Problem: Memory issues with 3D convolutions

```

1  # Use gradient checkpointing for large models
2  from torch.utils.checkpoint import checkpoint
3
4  class MemoryEfficientCNN(nn.Module):

```

```
5     def __init__(self):
6         super().__init__()
7         self.conv1 = nn.Conv3d(1, 32, 3)
8         self.conv2 = nn.Conv3d(32, 64, 3)
9         # ... more layers
10
11     def forward(self, x):
12         # Checkpoint expensive layers
13         x = checkpoint(self.conv1, x)
14         x = checkpoint(self.conv2, x)
15         return x
```

Tips for reducing memory:

- Use smaller batch sizes
- Use mixed precision training (FP16)
- Reduce number of filters
- Use gradient checkpointing
- Process data in patches (for very large images)

7.5 Exercises

7.1: Basic CNN for MNIST - ★★

Goal: Build and train your first CNN.

1. Build a CNN for MNIST: 2 conv layers + 2 FC layers
2. Use 32 and 64 filters, 3×3 kernels, ReLU, max pooling
3. Train for 5 epochs
4. Achieve >95% test accuracy
5. Visualize some predictions

Starter code:

```

1 from torchvision import datasets, transforms
2
3 # Load MNIST
4 transform = transforms.ToTensor()
5 train_dataset = datasets.MNIST('./data', train=True,
6                               download=True, transform=
7                               transform)
8 test_dataset = datasets.MNIST('./data', train=False,
9                               transform=transform)
10
11 # Your CNN here
12 class MyCNN(nn.Module):
13     def __init__(self):
14         super().__init__()
15         # Your architecture
16
17     def forward(self, x):
18         # Your forward pass
19         pass

```

7.2: Shape Calculation Practice - ★★

Goal: Master output size calculations.

For a 32×32 input image, calculate output sizes for:

1. Conv(3×3, stride=1, padding=1)
2. Conv(5×5, stride=2, padding=2)
3. Conv(3×3, stride=1, padding=0) followed by MaxPool(2×2)
4. Three consecutive Conv(3×3, stride=1, padding=1) layers
5. Conv(3×3, stride=1, padding=1, dilation=2)

Then implement these layers and verify your calculations with actual tensors.

7.3: 1D CNN for Time Series - ★★★

Goal: Apply convolutions to sequence data.

1. Generate synthetic time series: $y = \sin(x) + 0.5 \sin(3x) + \text{noise}$
2. Create dataset: windows of 100 points predict next 10 points
3. Build a 1D CNN: 3 conv layers with increasing filters
4. Train and compare with a simple MLP baseline
5. Visualize predictions on test data

Hint: Use `nn.Conv1d` with `kernel_size=7` or `9`.

7.4: VGG-Style Network - ★★★**Goal:** Build a deeper network with proper architecture.

1. Implement a VGG-style network:
 - Block 1: $2 \times \text{Conv}(64) + \text{Pool}$
 - Block 2: $2 \times \text{Conv}(128) + \text{Pool}$
 - Block 3: $3 \times \text{Conv}(256) + \text{Pool}$
 - Classifier: Global pooling + FC layers
2. Add batch normalization after each conv
3. Train on CIFAR-10
4. Achieve $>70\%$ test accuracy

Challenge: Add dropout and compare with/without it.**7.5: Simple U-Net - ★★★★★****Goal:** Implement encoder-decoder architecture.

Build a simple U-Net for image denoising:

1. Encoder: 3 conv blocks with downsampling
2. Decoder: 3 transposed conv blocks with upsampling
3. Add skip connections (concatenate encoder features with decoder)
4. Train on MNIST with added Gaussian noise
5. Visualize: noisy input \rightarrow denoised output \rightarrow clean target

Architecture:

Encoder:		Decoder:
1 \rightarrow 32	----->	concat \rightarrow 32 \rightarrow 1
32 \rightarrow 64	----->	concat \rightarrow 64 \rightarrow 32
64 \rightarrow 128	----->	128 \rightarrow 64

Starter code:

```

1 class SimpleUNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         # Encoder
5         self.enc1 = self.conv_block(1, 32)
6         self.enc2 = self.conv_block(32, 64)
7         self.enc3 = self.conv_block(64, 128)
8
9         # Decoder
10        self.dec3 = self.upconv_block(128, 64)
11        self.dec2 = self.upconv_block(128, 32) # 128 =
64 + 64 from skip
12        self.dec1 = nn.Conv2d(64, 1, 1) # 64 = 32 + 32
from skip
13
14        self.pool = nn.MaxPool2d(2)
15
16        def conv_block(self, in_ch, out_ch):
17            return nn.Sequential(
18                nn.Conv2d(in_ch, out_ch, 3, padding=1),
19                nn.BatchNorm2d(out_ch),
20                nn.ReLU()
21            )
22

```

```

23     def upconv_block(self, in_ch, out_ch):
24         return nn.Sequential(
25             nn.ConvTranspose2d(in_ch, out_ch, 2, stride
=2),
26             nn.BatchNorm2d(out_ch),
27             nn.ReLU()
28         )
29
30     def forward(self, x):
31         # Encoder with skip connections
32         e1 = self.enc1(x)
33         e2 = self.enc2(self.pool(e1))
34         e3 = self.enc3(self.pool(e2))
35
36         # Decoder with skip connections
37         d3 = self.dec3(e3)
38         d2 = self.dec2(torch.cat([d3, e2], dim=1)) #
Skip connection
39         d1 = self.dec1(torch.cat([d2, e1], dim=1)) #
Skip connection
40
41         return torch.sigmoid(d1)

```

7.6: 3D CNN for Volumetric Data - ★★★★★

Goal: Work with 3D convolutions.

1. Generate synthetic 3D data: 3D Gaussian blobs as "tumors"
2. Build a 3D CNN classifier: 3 conv layers + global pooling + FC
3. Handle memory constraints (small batch size, fewer filters)
4. Achieve good classification accuracy
5. Visualize 3D volumes and predictions (use slicing)

Data generation:

```

1 def generate_3d_blob(size=32):
2     """Generate 3D volume with central blob."""
3     volume = np.zeros((size, size, size))
4     center = size // 2
5
6     # Create Gaussian blob
7     for i in range(size):
8         for j in range(size):
9             for k in range(size):
10                 dist = np.sqrt((i-center)**2 + (j-center
)**2 +
11                               (k-center)**2)
12                 volume[i, j, k] = np.exp(-dist**2 / 100)
13
14     return torch.FloatTensor(volume).unsqueeze(0) # Add
channel dim

```

Warning: Use batch_size=2 or 4 for 3D convolutions!

7.6 Key Takeaways

Why CNNs work:

- **Parameter sharing:** Drastically fewer parameters than MLPs
- **Translation invariance:** Features work regardless of position
- **Hierarchical learning:** Low-level \rightarrow high-level features
- **Local connectivity:** Exploits spatial structure

Architecture design principles:

- Start with 32-64 filters, double after each pooling
- Use small kernels (3×3) rather than large ones
- Add batch normalization for deep networks
- Use strided convolutions or pooling for downsampling
- Global pooling at the end for flexibility

Convolution variants:

- **1D:** Time series, audio, text sequences
- **2D:** Images (most common)
- **3D:** Videos, medical scans, volumetric data
- **Transposed:** Upsampling in decoder/generator
- **Dilated:** Large receptive field without losing resolution

Common patterns:

- **VGG-style:** Stack conv blocks, double channels after pooling
- **Encoder-decoder:** Downsample then upsample (autoencoders, U-Net)
- **Skip connections:** Connect encoder to decoder (U-Net, ResNet)

Debugging tips:

- Always check shapes at each layer
- Visualize feature maps to understand what network learns
- Start small (few layers, few filters) then scale up
- Use `AdaptiveAvgPool2d` for flexibility
- Watch out for memory with 3D convolutions

8 Residual Networks & Skip Connections

8.1 Introduction: The Deep Network Problem

In the early 2010s, researchers discovered something puzzling: making networks deeper didn't always make them better. Sometimes, a 56-layer network performed **worse** than a 20-layer network—even on training data!

This wasn't overfitting (test loss was also worse). It was a fundamental optimization problem.

The Degradation Problem:

- Very deep networks are harder to optimize
- Gradients vanish or explode
- Training error increases with depth (counterintuitively)
- Not caused by overfitting—happens on training set too

Residual Networks (ResNets) solved this in 2015, enabling networks with 100+ layers that actually train better than shallow ones. The key insight: **skip connections** (also called **residual connections**).

Impact:

- ResNet-152 won ImageNet 2015 (3.57% error, human-level performance)
- Now standard in almost all vision architectures
- Fundamental principle used across all deep learning (Transformers, etc.)

8.2 Theory: Why Skip Connections Work

8.2.1 The Problem: Vanishing Gradients

In a deep network without skip connections:

$$\mathbf{x}_L = f_L(f_{L-1}(\dots f_2(f_1(\mathbf{x}_0))))$$

During backpropagation, gradients multiply through the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_0} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_L} \cdot \frac{\partial \mathbf{x}_L}{\partial \mathbf{x}_{L-1}} \cdot \dots \cdot \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1} \cdot \frac{\partial \mathbf{x}_1}{\partial \mathbf{x}_0}$$

If each term is less than 1, the gradient **vanishes** exponentially as it propagates backward.

Example: If each layer multiplies gradient by 0.9:

- After 10 layers: $0.9^{10} \approx 0.35$
- After 20 layers: $0.9^{20} \approx 0.12$
- After 50 layers: $0.9^{50} \approx 0.005$ (essentially zero!)

Early layers barely learn, making deep networks ineffective.

8.2.2 The Solution: Residual Connections

Instead of learning a direct mapping $\mathbf{y} = \mathcal{F}(\mathbf{x})$, learn a **residual**:

$$\mathbf{y} = \mathcal{F}(\mathbf{x}) + \mathbf{x}$$

The network learns the **difference** (residual) from the identity mapping.

[Key Insight: Identity Mapping] With a skip connection, the network can learn the identity function by simply setting $\mathcal{F}(\mathbf{x}) = 0$ (all weights to zero).

Why this matters:

- Identity is easy to learn (do nothing)
- Network starts with identity, then learns refinements
- Even if \mathcal{F} is poorly initialized, gradients still flow through skip connection
- Optimization becomes: "How should I modify the input?" instead of "What should the output be?"

8.2.3 Gradient Flow Through Skip Connections

During backpropagation with skip connections:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} \left(\frac{\partial \mathcal{F}}{\partial \mathbf{x}} + 1 \right)$$

The $+1$ term creates a **gradient highway**:

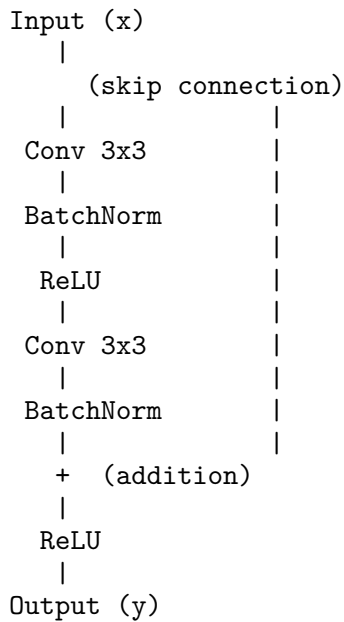
- Gradients always have a direct path backward (through the $+1$)
- Even if $\frac{\partial \mathcal{F}}{\partial \mathbf{x}} \approx 0$, gradients still flow
- No vanishing gradient problem!

Analogy: Regular network = narrow mountain trail (one path, easy to get stuck)

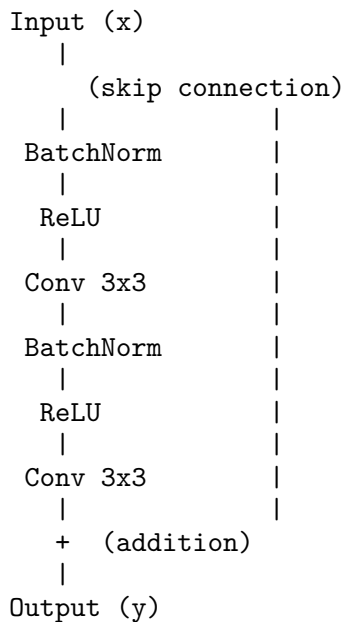
Skip connections = trail + highway (gradient can always get through)

8.2.4 Residual Block Variants

Basic Residual Block (Original):



Pre-Activation Residual Block (Improved):



Difference:

- **Post-activation:** Conv \rightarrow BN \rightarrow ReLU (then add)
- **Pre-activation:** BN \rightarrow ReLU \rightarrow Conv (then add)

Pre-activation advantages:

- Better gradient flow (identity mapping is truly unimpeded)
- Easier optimization
- Better for very deep networks (>100 layers)

8.2.5 When Dimensions Don't Match

If input and output have different dimensions, use a **projection shortcut**:

$$\mathbf{y} = \mathcal{F}(\mathbf{x}) + \mathbf{W}_s \mathbf{x}$$

where \mathbf{W}_s is a linear projection (1×1 convolution).

Two approaches:

1. **Identity shortcut + zero padding:** Pad channels with zeros (parameter-free)
2. **Projection shortcut:** Use 1×1 conv to match dimensions (adds parameters)

Most implementations use projection shortcuts when dimensions change.

8.3 Implementation: Building Residual Blocks

8.3.1 Basic Residual Block

```

1 import torch
2 import torch.nn as nn
3
4 class BasicBlock(nn.Module):
5     """Basic residual block with two 3x3 convolutions."""
6
7     def __init__(self, in_channels, out_channels, stride=1):
8         super().__init__()
9
10        # Main path
11        self.conv1 = nn.Conv2d(in_channels, out_channels,
12                                kernel_size=3, stride=stride,
13                                padding=1,
14                                bias=False)
15        self.bn1 = nn.BatchNorm2d(out_channels)
16        self.relu = nn.ReLU(inplace=True)
17
18        self.conv2 = nn.Conv2d(out_channels, out_channels,
19                                kernel_size=3, stride=1,
20                                padding=1,
21                                bias=False)
22        self.bn2 = nn.BatchNorm2d(out_channels)
23
24        # Skip connection (identity or projection)
25        self.skip = nn.Sequential()
26        if stride != 1 or in_channels != out_channels:
27            self.skip = nn.Sequential(
28                nn.Conv2d(in_channels, out_channels,
29                            kernel_size=1, stride=stride, bias=
30                            False),
31                nn.BatchNorm2d(out_channels)
32            )
33
34        def forward(self, x):
35            identity = x
36
37            # Main path
38            out = self.conv1(x)
39            out = self.bn1(out)
40            out = self.relu(out)
41
42            out = self.conv2(out)

```

```
40         out = self.bn2(out)
41
42         # Skip connection
43         identity = self.skip(identity)
44
45         # Add residual
46         out += identity
47         out = self.relu(out)
48
49         return out
50
51 # Test
52 block = BasicBlock(64, 64)
53 x = torch.randn(4, 64, 32, 32)
54 y = block(x)
55 print(y.shape) # torch.Size([4, 64, 32, 32])
56
57 # Test with dimension change
58 block_downsample = BasicBlock(64, 128, stride=2)
59 y2 = block_downsample(x)
60 print(y2.shape) # torch.Size([4, 128, 16, 16])
```

[Why bias=False?] When using batch normalization immediately after convolution, the bias is redundant (BatchNorm has its own bias term). Setting `bias=False` saves parameters and computation without affecting performance.

8.3.2 Pre-Activation Residual Block

```

1 class PreActBlock(nn.Module):
2     """Pre-activation residual block (improved version)."""
3
4     def __init__(self, in_channels, out_channels, stride=1):
5         super().__init__()
6
7         # Pre-activation
8         self.bn1 = nn.BatchNorm2d(in_channels)
9         self.relu = nn.ReLU(inplace=True)
10        self.conv1 = nn.Conv2d(in_channels, out_channels,
11                                kernel_size=3, stride=stride,
12                                padding=1,
13                                bias=False)
14
15        self.bn2 = nn.BatchNorm2d(out_channels)
16        self.conv2 = nn.Conv2d(out_channels, out_channels,
17                                kernel_size=3, stride=1,
18                                padding=1,
19                                bias=False)
20
21        # Skip connection
22        self.skip = nn.Sequential()
23        if stride != 1 or in_channels != out_channels:
24            self.skip = nn.Sequential(
25                nn.Conv2d(in_channels, out_channels,
26                           kernel_size=1, stride=stride, bias=
27                           False)
28            )
29
30    def forward(self, x):
31        # Pre-activation
32        out = self.bn1(x)
33        out = self.relu(out)
34
35        # Save for skip connection (after activation)
36        identity = self.skip(out)
37
38        # First conv
39        out = self.conv1(out)
40
41        # Second pre-activation + conv
42        out = self.bn2(out)
43        out = self.relu(out)
44        out = self.conv2(out)
45
46        # Add residual (no activation after addition!)
47        out += identity
48
49        return out

```

Key difference: In pre-activation, the final addition has no activation. The next block will apply activation first.

8.3.3 Bottleneck Block

For deeper networks (ResNet-50+), use bottleneck blocks to reduce parameters:

```

1 class BottleneckBlock(nn.Module):
2     """
3     Bottleneck block: 1x1 -> 3x3 -> 1x1

```

```

4  Reduces parameters by using 1x1 convolutions to reduce/
   expand channels.
5  """
6  expansion = 4  # Output channels = input channels * 4
7
8  def __init__(self, in_channels, out_channels, stride=1):
9      super().__init__()
10
11     # Bottleneck: reduce channels
12     self.conv1 = nn.Conv2d(in_channels, out_channels,
13                             kernel_size=1, bias=False)
14     self.bn1 = nn.BatchNorm2d(out_channels)
15
16     # Main 3x3 conv
17     self.conv2 = nn.Conv2d(out_channels, out_channels,
18                             kernel_size=3, stride=stride,
19                             padding=1, bias=False)
20     self.bn2 = nn.BatchNorm2d(out_channels)
21
22     # Expand channels
23     self.conv3 = nn.Conv2d(out_channels, out_channels *
self.expansion,
24                             kernel_size=1, bias=False)
25     self.bn3 = nn.BatchNorm2d(out_channels * self.
expansion)
26
27     self.relu = nn.ReLU(inplace=True)
28
29     # Skip connection
30     self.skip = nn.Sequential()
31     if stride != 1 or in_channels != out_channels * self.
expansion:
32         self.skip = nn.Sequential(
33             nn.Conv2d(in_channels, out_channels * self.
expansion,
34                         kernel_size=1, stride=stride, bias=
False),
35             nn.BatchNorm2d(out_channels * self.expansion)
36         )
37
38     def forward(self, x):
39         identity = x
40
41         out = self.conv1(x)
42         out = self.bn1(out)
43         out = self.relu(out)
44
45         out = self.conv2(out)
46         out = self.bn2(out)
47         out = self.relu(out)
48
49         out = self.conv3(out)
50         out = self.bn3(out)
51
52         identity = self.skip(identity)
53         out += identity
54         out = self.relu(out)
55
56         return out
57
58 # Compare parameters
59 basic = BasicBlock(256, 256)

```

```
60 bottleneck = BottleneckBlock(256, 64) # Output: 64*4 = 256
61
62 basic_params = sum(p.numel() for p in basic.parameters())
63 bottleneck_params = sum(p.numel() for p in bottleneck.
    parameters())
64
65 print(f"Basic block: {basic_params:,} parameters")
66 print(f"Bottleneck: {bottleneck_params:,} parameters")
67 # Basic block: ~590,000 parameters
68 # Bottleneck: ~70,000 parameters (much fewer!)
```

8.3.4 Building a Complete ResNet

```

1 class ResNet(nn.Module):
2     """ResNet architecture for CIFAR-10."""
3
4     def __init__(self, block, num_blocks, num_classes=10):
5         """
6         Args:
7             block: BasicBlock or BottleneckBlock
8             num_blocks: List of number of blocks per layer
9             num_classes: Number of output classes
10        """
11        super().__init__()
12        self.in_channels = 64
13
14        # Initial convolution
15        self.conv1 = nn.Conv2d(3, 64, kernel_size=3,
16                                stride=1, padding=1, bias=False)
17
18        self.bn1 = nn.BatchNorm2d(64)
19        self.relu = nn.ReLU(inplace=True)
20
21        # Residual layers
22        self.layer1 = self._make_layer(block, 64, num_blocks
23[0], stride=1)
24        self.layer2 = self._make_layer(block, 128, num_blocks
25[1], stride=2)
26        self.layer3 = self._make_layer(block, 256, num_blocks
27[2], stride=2)
28        self.layer4 = self._make_layer(block, 512, num_blocks
29[3], stride=2)
30
31        # Global pooling and classifier
32        self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
33        self.fc = nn.Linear(512 * block.expansion,
34                                num_classes)
35
36        def _make_layer(self, block, out_channels, num_blocks,
37                        stride):
38            """Create a layer with multiple residual blocks."""
39            strides = [stride] + [1] * (num_blocks - 1)
40            layers = []
41
42            for stride in strides:
43                layers.append(block(self.in_channels,
44                                out_channels, stride))
45                self.in_channels = out_channels * block.expansion
46
47            return nn.Sequential(*layers)
48
49        def forward(self, x):
50            out = self.conv1(x)
51            out = self.bn1(out)
52            out = self.relu(out)
53
54            out = self.layer1(out)
55            out = self.layer2(out)
56            out = self.layer3(out)
57            out = self.layer4(out)
58
59            out = self.avgpool(out)
60            out = out.view(out.size(0), -1)

```

```

53         out = self.fc(out)
54
55         return out
56
57     # ResNet-18
58     def ResNet18(num_classes=10):
59         return ResNet(BasicBlock, [2, 2, 2, 2], num_classes)
60
61     # ResNet-34
62     def ResNet34(num_classes=10):
63         return ResNet(BasicBlock, [3, 4, 6, 3], num_classes)
64
65     # ResNet-50 (uses bottleneck)
66     def ResNet50(num_classes=10):
67         # Note: BottleneckBlock needs expansion=4
68         BottleneckBlock.expansion = 4
69         return ResNet(BottleneckBlock, [3, 4, 6, 3], num_classes)
70
71     # Test
72     model = ResNet18()
73     x = torch.randn(4, 3, 32, 32)
74     y = model(x)
75     print(y.shape) # torch.Size([4, 10])
76
77     # Count parameters
78     total_params = sum(p.numel() for p in model.parameters())
79     print(f"ResNet-18: {total_params:,} parameters")

```

[ResNet Naming] ResNet-N refers to the total number of weight layers:

ResNet-18: 1 initial conv + 4 layers \times 2 blocks \times 2 convs + 1 FC = 18 layers

ResNet-34: 1 + (3+4+6+3) \times 2 + 1 = 34 layers

ResNet-50: 1 + (3+4+6+3) \times 3 + 1 = 50 layers (bottleneck has 3 convs)

Deeper = more capacity but slower training/inference.

8.3.5 Dense Connections (DenseNet)

DenseNet takes skip connections further: **every layer connects to every other layer**.

```

1 class DenseBlock(nn.Module):
2     """Dense block where each layer receives all previous
3     layers as input."""
4
5     def __init__(self, in_channels, growth_rate, num_layers):
6         """
7         Args:
8             in_channels: Initial number of channels
9             growth_rate: Number of channels added per layer
10            num_layers: Number of layers in the dense block
11        """
12        super().__init__()
13
14        self.layers = nn.ModuleList()
15        for i in range(num_layers):
16            self.layers.append(
17                self._make_dense_layer(in_channels + i *
18                                     growth_rate,
19                                     growth_rate)
20            )
21
22        def _make_dense_layer(self, in_channels, growth_rate):
23            return nn.Sequential(
24                nn.BatchNorm2d(in_channels),
25                nn.ReLU(inplace=True),
26                nn.Conv2d(in_channels, growth_rate, kernel_size
27                          =3,
28                          padding=1, bias=False)
29            )
30
31        def forward(self, x):
32            features = [x]
33            for layer in self.layers:
34                # Concatenate all previous feature maps
35                new_features = layer(torch.cat(features, dim=1))
36                features.append(new_features)
37
38            # Return concatenation of all features
39            return torch.cat(features, dim=1)
40
41 # Test
42 block = DenseBlock(in_channels=64, growth_rate=32, num_layers
43                   =4)
44 x = torch.randn(2, 64, 16, 16)
45 y = block(x)
46 print(y.shape) # torch.Size([2, 192, 16, 16])
47 # 192 = 64 (input) + 4 * 32 (4 layers, each adds 32 channels)

```

DenseNet characteristics:

- **Feature reuse:** All layers have access to all previous features
- **Fewer parameters:** Smaller growth_rate than ResNet channel counts
- **Memory intensive:** Must store all intermediate features
- **Better gradient flow:** Even better than ResNet

When to use DenseNet vs ResNet:

- **DenseNet:** When parameter efficiency is critical, smaller datasets
- **ResNet:** When memory is limited, faster training/inference, more standard

8.3.6 Skip Connections in Practice

General principles for adding skip connections:

1. **Add skip every 2-3 layers:** Balance between gradient flow and complexity
2. **Match dimensions:** Use projection when channels or spatial size change
3. **Identity when possible:** Projection-free shortcuts are better
4. **Add before final activation:** Preserves identity mapping
5. **Don't overdo it:** Too many skip connections can hurt performance

Example: Adding skip connections to existing MLP:

```

1 class MLPWithSkipConnections(nn.Module):
2     """MLP with residual connections."""
3
4     def __init__(self, input_dim, hidden_dim, output_dim,
5 num_layers=4):
6         super().__init__()
7
8         self.input_layer = nn.Linear(input_dim, hidden_dim)
9
10        # Hidden layers with skip connections
11        self.hidden_layers = nn.ModuleList()
12        for _ in range(num_layers):
13            self.hidden_layers.append(
14                nn.Sequential(
15                    nn.Linear(hidden_dim, hidden_dim),
16                    nn.ReLU(),
17                    nn.Linear(hidden_dim, hidden_dim)
18                )
19            )
20
21        self.output_layer = nn.Linear(hidden_dim, output_dim)
22        self.relu = nn.ReLU()
23
24    def forward(self, x):
25        x = self.relu(self.input_layer(x))
26
27        # Apply each hidden layer with skip connection
28        for layer in self.hidden_layers:
29            residual = x
30            x = layer(x)
31            x = x + residual # Skip connection
32            x = self.relu(x)
33
34        x = self.output_layer(x)
35        return x

```

8.4 Implementation: Debugging ResNets

8.4.1 Verifying Skip Connections

```

1 def verify_skip_connections(model, input_shape=(1, 3, 32, 32)
  ):
2     """Verify that skip connections are working."""
3     x = torch.randn(*input_shape, requires_grad=True)
4
5     # Forward pass
6     output = model(x)
7     loss = output.sum()
8
9     # Backward pass
10    loss.backward()
11
12    # Check if gradients reach input
13    if x.grad is not None:
14        print(f"    Gradients flow to input: {x.grad.abs().
mean().item():.6f}")
15    else:
16        print("    No gradients at input!")
17
18    # Check gradient magnitudes throughout network
19    print("\nGradient magnitudes by layer:")
20    for name, param in model.named_parameters():
21        if param.grad is not None:
22            grad_mean = param.grad.abs().mean().item()
23            grad_max = param.grad.abs().max().item()
24            print(f"{name:30s}: mean={grad_mean:.6f}, max={
grad_max:.6f}")
25
26    # Usage
27    # model = ResNet18()
28    # verify_skip_connections(model)

```

8.4.2 Comparing With and Without Skip Connections

```

1 class PlainCNN(nn.Module):
2     """Same architecture as ResNet but without skip
connections."""
3
4     def __init__(self, num_classes=10):
5         super().__init__()
6
7         self.conv1 = nn.Conv2d(3, 64, 3, padding=1, bias=
False)
8         self.bn1 = nn.BatchNorm2d(64)
9
10        # Same structure as ResNet-18 but no skip connections
11        self.layers = nn.Sequential(
12            # Layer 1
13            nn.Conv2d(64, 64, 3, padding=1, bias=False),
14            nn.BatchNorm2d(64),
15            nn.ReLU(),
16            nn.Conv2d(64, 64, 3, padding=1, bias=False),
17            nn.BatchNorm2d(64),
18            nn.ReLU(),
19
20            # Layer 2 (with stride=2 downsampling)
21            nn.Conv2d(64, 128, 3, stride=2, padding=1, bias=
False),

```

```

22         nn.BatchNorm2d(128),
23         nn.ReLU(),
24         nn.Conv2d(128, 128, 3, padding=1, bias=False),
25         nn.BatchNorm2d(128),
26         nn.ReLU(),
27
28         # Continue pattern...
29     )
30
31     self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
32     self.fc = nn.Linear(512, num_classes)
33
34     def forward(self, x):
35         x = nn.ReLU()(self.bn1(self.conv1(x)))
36         x = self.layers(x)
37         x = self.avgpool(x)
38         x = x.view(x.size(0), -1)
39         return self.fc(x)
40
41 # Compare training curves
42 # plain_model = PlainCNN()
43 # resnet_model = ResNet18()
44 # Train both and plot losses - ResNet should converge faster
   and better

```

8.5 Exercises

8.1: Basic Residual Block - ★★

Goal: Implement and understand a basic residual block.

1. Implement a `BasicBlock` with two 3×3 convolutions
2. Include skip connection with projection for dimension changes
3. Test with different input/output channels:
 - Same channels ($64 \rightarrow 64$)
 - Different channels ($64 \rightarrow 128$)
 - With downsampling ($\text{stride}=2$)
4. Verify output shapes match expected dimensions

Starter code:

```

1 class MyBasicBlock(nn.Module):
2     def __init__(self, in_channels, out_channels, stride
   =1):
3         super().__init__()
4         # Your code here
5
6     def forward(self, x):
7         # Your code here
8         pass
9
10 # Test cases
11 block1 = MyBasicBlock(64, 64, stride=1)
12 block2 = MyBasicBlock(64, 128, stride=2)
13
14 x = torch.randn(4, 64, 32, 32)
15 print(block1(x).shape) # Should be [4, 64, 32, 32]
16 print(block2(x).shape) # Should be [4, 128, 16, 16]

```

8.2: Plain vs Residual Comparison - ★★★**Goal:** Empirically verify that skip connections help.

1. Build two identical networks (10 layers each):
 - One **without** skip connections (plain CNN)
 - One **with** skip connections (ResNet-style)
2. Train both on CIFAR-10 or MNIST
3. Track and plot:
 - Training loss
 - Validation accuracy
 - Gradient magnitudes in early layers
4. Observe: ResNet should train faster and achieve better accuracy

Questions to answer:

- How much faster does ResNet converge?
- What's the final accuracy difference?
- How do gradient magnitudes compare?

8.3: Pre-Activation vs Post-Activation - ★★★**Goal:** Compare residual block variants.

1. Implement both post-activation and pre-activation blocks
2. Build identical ResNets with each type
3. Train on the same dataset
4. Compare:
 - Training stability (loss curves)
 - Final accuracy
 - Training time

Hypothesis: Pre-activation should be slightly better, especially for deeper networks.**Extension:** Try with different depths (18, 34, 50 layers) and see if the gap widens.**8.4: Building ResNet-18 - ★★★****Goal:** Build a complete ResNet from scratch.

1. Implement ResNet-18 architecture:
 - Initial 7×7 conv (or 3×3 for CIFAR)
 - 4 residual layers with $[2, 2, 2, 2]$ blocks
 - Channels: $64 \rightarrow 128 \rightarrow 256 \rightarrow 512$
 - Global average pooling + FC
2. Train on CIFAR-10
3. Achieve $>90\%$ test accuracy
4. Save the model and load it for inference

Bonus: Visualize learned features in first convolutional layer.**8.5: Bottleneck Block - ★★★★★****Goal:** Implement parameter-efficient bottleneck blocks.

1. Implement `BottleneckBlock`: $1 \times 1 \rightarrow 3 \times 3 \rightarrow 1 \times 1$

2. Compare parameter count with `BasicBlock` for same input/output
3. Build ResNet-50 using bottleneck blocks
4. Train and compare with ResNet-18:
 - Accuracy (ResNet-50 should be better)
 - Training time (ResNet-50 slower)
 - Memory usage

Challenge: How deep can you go before hitting memory limits? Try ResNet-101.

8.6: Skip Connections in U-Net - ★★★★★

Goal: Apply skip connections to encoder-decoder architectures.

1. Take the U-Net from CNN section (Exercise 7.5)
2. Add residual blocks in encoder and decoder paths
3. Train on image denoising task
4. Compare with original U-Net:
 - Denoising quality (PSNR/SSIM)
 - Training stability
 - Convergence speed

Architecture:

- Replace conv blocks with residual blocks
- Keep skip connections between encoder/decoder
- Use both residual blocks AND U-Net skip connections

8.7: Dense Block Implementation - ★★★★★

Goal: Implement and understand DenseNet.

1. Implement a `DenseBlock` where each layer receives all previous layers
2. Implement a transition layer (1×1 conv + pooling) to reduce channels
3. Build a small DenseNet:
 - 3 dense blocks with `growth_rate=12`
 - Transition layers between blocks
 - Final global pooling + FC
4. Compare with ResNet of similar depth:
 - Parameter count (DenseNet should be fewer)
 - Memory usage (DenseNet higher during training)
 - Accuracy

Hint: Use `torch.cat()` to concatenate feature maps along channel dimension.

8.6 Key Takeaways

The fundamental problem:

- Very deep networks suffer from degradation (not just overfitting)
- Vanishing gradients make early layers hard to train
- Optimization becomes difficult, not just capacity

How skip connections solve it:

- Create gradient highways: $\frac{\partial \mathcal{L}}{\partial \mathbf{x}}$ always has +1 term
- Learn residuals instead of direct mappings (easier optimization)
- Enable training of 100+ layer networks
- Identity mapping is easy to learn (set weights to zero)

Residual block variants:

- **Post-activation:** Conv \rightarrow BN \rightarrow ReLU, then add (original)
- **Pre-activation:** BN \rightarrow ReLU \rightarrow Conv, then add (improved)
- **Bottleneck:** $1 \times 1 \rightarrow 3 \times 3 \rightarrow 1 \times 1$ (for deeper networks, fewer parameters)
- **Dense:** Every layer connects to all previous layers (extreme skip connections)

Implementation guidelines:

- Use projection shortcuts when dimensions change
- Set `bias=False` when using BatchNorm
- Add skip connections every 2-3 layers
- Pre-activation is better for very deep networks (>100 layers)
- Use bottleneck blocks for ResNet-50+

When to use skip connections:

- **Always use for deep networks (>20 layers)**
- **Vision tasks:** ResNet is standard backbone
- **U-Net style architectures:** Encoder-decoder with skip connections
- **Any architecture where gradient flow is critical**
- **When training is unstable or slow**

Architecture choices:

Model	Layers	Use Case	Parameters
ResNet-18	18	General purpose, fast	11M
ResNet-34	34	More capacity	21M
ResNet-50	50	High accuracy	25M
ResNet-101	101	Maximum accuracy	44M
DenseNet-121	121	Parameter efficient	8M

Table 5: Popular ResNet/DenseNet architectures

Impact on deep learning:

- Enabled training of 100+ layer networks
- Now standard in almost all vision models

- Extended to other domains (Transformers use similar ideas)
- Fundamental principle: make it easy for network to learn identity

Debugging tips:

- Verify gradients flow to early layers
- Compare with plain network (without skip connections)
- Check gradient magnitudes throughout network
- Visualize activations in residual blocks
- If training fails, check skip connection implementation

Common mistakes:

1. Forgetting projection when dimensions change
2. Applying activation after addition in pre-activation blocks
3. Not using BatchNorm (skip connections alone aren't enough)
4. Making skip connection too complex (keep it identity when possible)
5. Adding too many skip connections (can hurt performance)

[Rule of Thumb] **For networks >20 layers:** Always use residual connections

For networks >50 layers: Use bottleneck blocks

For networks >100 layers: Use pre-activation blocks

When in doubt: Use ResNet as your starting point and modify from there.

9 Batch Normalization & Layer Normalization

9.1 Introduction: Why Normalization Matters

Training deep networks is hard. As gradients flow backward, they can vanish or explode. As data flows forward, the distribution of activations can shift. These issues slow training and make networks sensitive to initialization.

Normalization techniques stabilize training by controlling the distribution of activations.

Key benefits:

- **Faster training:** Can use higher learning rates (2-10× speedup)
- **Less sensitivity to initialization:** Network more robust
- **Regularization effect:** Acts like dropout (slight noise from batch statistics)
- **Enables deeper networks:** Makes 100+ layer networks trainable

When you need normalization:

- Deep networks (>10 layers)
- Training is slow or unstable
- Network sensitive to learning rate or initialization
- Using high learning rates

9.2 Theory: Batch Normalization

9.2.1 The Problem: Internal Covariate Shift

As we train, the distribution of inputs to each layer changes. This is called **internal covariate shift**.

Example: Consider layer 3 of a network.

- Early in training: inputs to layer 3 might have mean=0.5, std=0.3
- After 100 steps: mean=2.1, std=1.8 (distribution has shifted!)

Each layer must constantly adapt to the changing distribution. This slows learning.

Batch Normalization's solution: Normalize inputs to each layer to have consistent statistics.

9.2.2 How Batch Normalization Works

Given a batch of inputs $\mathbf{x} = \{x_1, x_2, \dots, x_m\}$:

Step 1: Compute batch statistics

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i \quad (\text{batch mean})$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad (\text{batch variance})$$

Step 2: Normalize

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

where ϵ is a small constant (typically 10^{-5}) for numerical stability.

Step 3: Scale and shift (learnable parameters)

$$y_i = \gamma \hat{x}_i + \beta$$

where γ (scale) and β (shift) are learned parameters.

Why scale and shift? The network can learn to undo the normalization if needed! If $\gamma = \sqrt{\sigma_B^2}$ and $\beta = \mu_B$, we recover the original distribution.

[Key Insight] Batch norm has two modes:

Training: Use batch statistics (μ_B, σ_B^2)

- Computed from current batch
- Adds stochasticity (different batches \rightarrow different statistics)
- Acts as regularization

Inference: Use running statistics ($\mu_{running}, \sigma_{running}^2$)

- Exponential moving average from training
- Deterministic (same input \rightarrow same output)
- No batch dependency

During training, PyTorch maintains:

$$\mu_{running} \leftarrow (1 - \text{momentum}) \cdot \mu_{running} + \text{momentum} \cdot \mu_B$$

Default momentum = 0.1.

9.2.3 Why Batch Normalization Helps

1. Smooths the optimization landscape

Batch norm makes the loss surface smoother, allowing larger learning rates.

2. Reduces internal covariate shift

Each layer receives inputs with consistent statistics.

3. Acts as regularization

Noise from batch statistics acts like dropout (but usually weaker).

4. Reduces sensitivity to initialization

Even with poor initialization, batch norm helps normalize activations.

5. Allows higher learning rates

The smoothing effect permits learning rates 10× higher than without batch norm.

9.2.4 Limitations of Batch Normalization

1. Requires large enough batches

With `batch_size < 8`, batch statistics are unreliable. Solution: Use `LayerNorm` or `GroupNorm`.

2. Different behavior in train/eval

Must remember to call `model.train()` and `model.eval()`. Forgetting this is a common bug!

3. Doesn't work well for RNNs

Batch statistics across sequence positions don't make sense. Solution: Use `LayerNorm`.

4. Coupling between samples in batch

Each sample's normalization depends on other samples in the batch. Can cause issues in some scenarios.

9.3 Theory: Layer Normalization

Layer Normalization normalizes across features instead of across the batch.

Batch Norm: Normalize each feature across the batch

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad \text{where } \mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij}$$

(Averages over batch dimension i for each feature j)

Layer Norm: Normalize all features for each sample

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} \quad \text{where } \mu_i = \frac{1}{d} \sum_{j=1}^d x_{ij}$$

(Averages over feature dimension j for each sample i)

Key differences:

Aspect	Batch Norm	Layer Norm
Normalizes over	Batch dimension	Feature dimension
Batch size dependency	Yes (fails for small batches)	No (works with batch=1)
Train/eval modes	Different	Same
Learnable params	γ, β per feature	γ, β per feature
Best for	CNNs, MLPs	RNNs, Transformers
Running statistics	Yes	No

Table 6: Batch Norm vs Layer Norm**When to use Layer Norm:**

- RNNs and sequence models
- Transformers (standard choice)
- Small batch sizes (`batch_size < 8`)
- Online learning (single sample at a time)
- When you want same behavior in train/eval

9.4 Implementation: Using Normalization in PyTorch

9.4.1 Batch Normalization Usage

```

1 import torch
2 import torch.nn as nn
3
4 # For 1D data (MLPs, sequences)
5 bn1d = nn.BatchNorm1d(num_features=128)
6
7 # For 2D data (CNNs)
8 bn2d = nn.BatchNorm2d(num_features=64)
9
10 # For 3D data (3D CNNs, video)
11 bn3d = nn.BatchNorm3d(num_features=32)
12
13 # Example: CNN with batch norm
14 model = nn.Sequential(
15     nn.Conv2d(3, 64, 3, padding=1),
16     nn.BatchNorm2d(64), # Batch norm after conv
17     nn.ReLU(),
18
19     nn.Conv2d(64, 128, 3, padding=1),
20     nn.BatchNorm2d(128),
21     nn.ReLU(),
22 )
23
24 # Training mode (default)
25 model.train()
26 x = torch.randn(16, 3, 32, 32)
27 out = model(x) # Uses batch statistics
28
29 # Evaluation mode
30 model.eval()
31 with torch.no_grad():
32     out = model(x) # Uses running statistics

```

Where to place batch norm:

Option 1: Conv → BN → Activation (recommended)

```

1 nn.Conv2d(in_ch, out_ch, 3, padding=1),
2 nn.BatchNorm2d(out_ch),
3 nn.ReLU()

```

Option 2: Conv → Activation → BN (less common)

```

1 nn.Conv2d(in_ch, out_ch, 3, padding=1),
2 nn.ReLU(),
3 nn.BatchNorm2d(out_ch)

```

Both work, but Option 1 is more standard.

[Critical: Set bias=False] When using batch norm after a linear/conv layer, set bias=False:

```

1 # WRONG: Batch norm makes bias redundant
2 nn.Conv2d(3, 64, 3, padding=1, bias=True), # Wasteful!
3 nn.BatchNorm2d(64), # Has its own bias (beta)
4
5 # CORRECT: No need for conv bias
6 nn.Conv2d(3, 64, 3, padding=1, bias=False), # No bias

```

```
7 nn.BatchNorm2d(64), # Beta parameter serves as bias
```

Why? Batch norm subtracts the mean, so any constant bias gets canceled out. The β parameter in batch norm serves as the bias.

9.4.2 Layer Normalization Usage

```
1 # Layer norm normalizes over feature dimension
2 ln = nn.LayerNorm(normalized_shape=128)
3
4 # For sequences: (batch, seq_len, features)
5 x = torch.randn(32, 100, 128) # 32 sequences, length 100,
   128 features
6 out = ln(x) # Normalizes over the 128 features for each
   position
7
8 # Example: Transformer-style layer
9 class TransformerBlock(nn.Module):
10     def __init__(self, d_model):
11         super().__init__()
12         self.attention = nn.MultiheadAttention(d_model,
   num_heads=8)
13         self.ln1 = nn.LayerNorm(d_model)
14         self.ffn = nn.Sequential(
15             nn.Linear(d_model, 4*d_model),
16             nn.ReLU(),
17             nn.Linear(4*d_model, d_model)
18         )
19         self.ln2 = nn.LayerNorm(d_model)
20
21     def forward(self, x):
22         # Attention with layer norm
23         attn_out, _ = self.attention(x, x, x)
24         x = self.ln1(x + attn_out) # Residual + norm
25
26         # FFN with layer norm
27         ffn_out = self.ffn(x)
28         x = self.ln2(x + ffn_out) # Residual + norm
29
30     return x
```

9.4.3 Other Normalization Variants

Instance Normalization (for style transfer):

```
1 # Normalizes each sample and channel independently
2 # Used in style transfer and GANs
3 in_norm = nn.InstanceNorm2d(num_features=64)
4
5 # Shape: (batch, channels, height, width)
6 # Computes mean/std for each (sample, channel) pair over (H, W)
```

Group Normalization (hybrid approach):

```
1 # Divides channels into groups, normalizes within each group
2 # Works well with small batch sizes
3 gn = nn.GroupNorm(num_groups=8, num_channels=64)
4
5 # Example: 64 channels      8 groups of 8 channels each
6 # Normalizes over spatial dimensions and within each group
```

When to use each:

Normalization	Use Case
Batch Norm	CNNs, MLPs, large batches (≥ 8)
Layer Norm	RNNs, Transformers, any batch size
Instance Norm	Style transfer, GANs
Group Norm	Small batch sizes, object detection

Table 7: Choosing normalization type

9.4.4 Implementing Batch Norm from Scratch

```
1 class MyBatchNorm1d(nn.Module):
2     """Batch normalization for understanding."""
3
4     def __init__(self, num_features, eps=1e-5, momentum=0.1):
5         super().__init__()
6
7         # Learnable parameters
8         self.gamma = nn.Parameter(torch.ones(num_features))
9         self.beta = nn.Parameter(torch.zeros(num_features))
10
11        # Running statistics (not trainable)
12        self.register_buffer('running_mean', torch.zeros(
num_features))
13        self.register_buffer('running_var', torch.ones(
num_features))
14
15        self.eps = eps
16        self.momentum = momentum
17
18        def forward(self, x):
19            # x shape: (batch, features)
20
21            if self.training:
22                # Training mode: use batch statistics
23                batch_mean = x.mean(dim=0)
```

```

24         batch_var = x.var(dim=0, unbiased=False)
25
26         # Normalize
27         x_norm = (x - batch_mean) / torch.sqrt(batch_var
+ self.eps)
28
29         # Update running statistics (exponential moving
average)
30         with torch.no_grad():
31             self.running_mean = (1 - self.momentum) *
self.running_mean + \
32                                     self.momentum * batch_mean
33             self.running_var = (1 - self.momentum) * self
.running_var + \
34                                     self.momentum * batch_var
35         else:
36             # Eval mode: use running statistics
37             x_norm = (x - self.running_mean) / \
38                     torch.sqrt(self.running_var + self.eps)
39
40         # Scale and shift
41         out = self.gamma * x_norm + self.beta
42
43         return out
44
45     # Test
46     bn = MyBatchNorm1d(10)
47     x = torch.randn(32, 10)
48
49     # Training mode
50     bn.train()
51     out_train = bn(x)
52
53     # Eval mode
54     bn.eval()
55     out_eval = bn(x)
56
57     print(f"Train output: {out_train.mean():.4f}, {out_train.std()
:.4f}")
58     print(f"Eval output: {out_eval.mean():.4f}, {out_eval.std()
:.4f}")

```


9.4.5 Common Mistakes and Debugging

Mistake 1: Forgetting to set eval mode

```

1 # WRONG: Model still in training mode during evaluation
2 model.train() # Set at beginning of training
3 # ... training loop ...
4 # Evaluation (FORGOT model.eval()!)
5 with torch.no_grad():
6     val_loss = evaluate(model, val_loader) # Uses batch
       statistics!
7
8 # CORRECT:
9 model.train()
10 # ... training ...
11 model.eval() # Switch to eval mode
12 with torch.no_grad():
13     val_loss = evaluate(model, val_loader) # Uses running
       statistics

```

Mistake 2: Small batch sizes with Batch Norm

```

1 # With batch_size=2, batch statistics are unreliable
2 # Symptoms: High variance in training, poor performance
3
4 # Solution 1: Increase batch size
5 train_loader = DataLoader(dataset, batch_size=32) # Instead
       of 2
6
7 # Solution 2: Use Layer Norm or Group Norm instead
8 # Replace BatchNorm2d with GroupNorm
9 nn.GroupNorm(num_groups=8, num_channels=64)

```

Mistake 3: Not setting bias=False

```

1 # WASTEFUL: Both conv and batch norm have biases
2 layer = nn.Sequential(
3     nn.Conv2d(64, 128, 3, bias=True), # Has bias
4     nn.BatchNorm2d(128) # Also has bias (beta)
5 )
6
7 # EFFICIENT: Only batch norm has bias
8 layer = nn.Sequential(
9     nn.Conv2d(64, 128, 3, bias=False), # No bias
10    nn.BatchNorm2d(128) # Beta serves as bias
11 )

```

Debugging: Check if batch norm is working

```

1 def check_batch_norm_stats(model, dataloader):
2     """Verify batch norm running statistics are being updated"""
3
4     model.train()
5
6     # Save initial running mean
7     bn_layer = None
8     for module in model.modules():
9         if isinstance(module, nn.BatchNorm2d):
10             bn_layer = module
11             break
12

```

```
13     if bn_layer is None:
14         print("No BatchNorm found!")
15         return
16
17     initial_mean = bn_layer.running_mean.clone()
18
19     # Run one batch
20     x, _ = next(iter(dataloader))
21     _ = model(x)
22
23     # Check if running mean changed
24     mean_changed = not torch.allclose(initial_mean, bn_layer.
25     running_mean)
26     print(f"Running mean updated: {mean_changed}")
27
28     # Check train vs eval difference
29     model.eval()
30     with torch.no_grad():
31         out_eval = model(x)
32
33     model.train()
34     out_train = model(x)
35
36     different = not torch.allclose(out_train, out_eval)
37     print(f"Train vs eval outputs differ: {different}")
38
39     # Usage
40     # check_batch_norm_stats(model, train_loader)
```

9.5 Exercises

9.1: Batch Norm Impact - ★★

Goal: See batch norm's effect empirically.

Train two identical networks on CIFAR-10:

1. **Without batch norm:** Plain CNN
2. **With batch norm:** Same CNN + batch norm after each conv

Compare:

- Training speed (epochs to reach 70% accuracy)
- Learning rate stability (try lr=0.01 and lr=0.1)
- Final test accuracy

Expected observation: Batch norm allows higher learning rate and converges faster.

Starter code:

```

1 class CNNWithoutBN(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.features = nn.Sequential(
5             nn.Conv2d(3, 64, 3, padding=1),
6             nn.ReLU(),
7             nn.Conv2d(64, 128, 3, padding=1),
8             nn.ReLU(),
9             nn.MaxPool2d(2),
10            # ... more layers
11        )
12
13 class CNNWithBN(nn.Module):
14     def __init__(self):
15         super().__init__()
16         self.features = nn.Sequential(
17             nn.Conv2d(3, 64, 3, padding=1, bias=False),
18             nn.BatchNorm2d(64),
19             nn.ReLU(),
20             # ... add batch norm after each conv
21        )

```

9.2: Train vs Eval Mode - ★★

Goal: Understand the importance of eval mode.

1. Train a CNN with batch norm on MNIST
2. After training, evaluate on test set **without** calling `model.eval()`
3. Evaluate again **with** `model.eval()`
4. Compare test accuracies
5. Investigate: Try with different batch sizes during evaluation (1, 8, 64)

Questions:

- How much does forgetting `eval()` hurt performance?
- Does batch size during eval matter when in train mode?
- Why does this happen?

9.3: Batch Norm from Scratch - ★★★**Goal:** Implement batch norm to understand it deeply.

1. Complete the `MyBatchNorm1d` implementation from above
2. Extend it to `MyBatchNorm2d` for CNNs
3. Test that it matches PyTorch's `nn.BatchNorm2d`:
 - Initialize both with same parameters
 - Feed same input
 - Verify outputs are close (use `torch.allclose`)
4. Train a small CNN using your custom batch norm
5. Verify training works correctly

Hint for BatchNorm2d:

```

1 # For input shape (batch, channels, height, width)
2 # Compute mean and var over dimensions (0, 2, 3)
3 # This gives per-channel statistics
4 batch_mean = x.mean(dim=(0, 2, 3), keepdim=True)
5 batch_var = x.var(dim=(0, 2, 3), keepdim=True, unbiased=False)

```

9.4: Layer Norm vs Batch Norm - ★★★**Goal:** Compare normalization types for sequences.

1. Generate a sequence classification task (e.g., classify sine vs cosine)
2. Build an RNN with:
 - Batch normalization
 - Layer normalization
 - No normalization
3. Train all three and compare:
 - Training stability
 - Final accuracy
 - Effect of batch size (try 4, 16, 64)

Expected result: Layer norm should work better for sequences, especially with small batches.**9.5: Placement Experiments - ★★★★★****Goal:** Investigate where to place batch norm.

Test different placements in a residual block:

1. Conv \rightarrow BN \rightarrow ReLU
2. Conv \rightarrow ReLU \rightarrow BN
3. BN \rightarrow ReLU \rightarrow Conv (pre-activation)

For each:

- Train on CIFAR-10
- Measure convergence speed
- Check gradient flow (from Section 8)
- Compare final accuracy

Questions:

- Which placement works best?

- Does the answer change with network depth?
- How does it affect gradient flow?

9.6: Small Batch Challenge - ★★★★★

Goal: Handle small batch sizes effectively.

1. Train a network with `batch_size=2` (deliberately small)
2. Try:
 - Batch Norm (observe instability)
 - Layer Norm
 - Group Norm (8 groups)
 - No normalization
3. Compare training stability and final performance
4. Investigate: Why does batch norm fail? Look at batch statistics variance.

Analysis code:

```

1  # Check batch statistics variance
2  bn_layer = model.features[1] # Assuming second layer is
   BN
3  means = []
4  vars = []
5
6  for batch in train_loader:
7      x, _ = batch
8      # Hook to capture batch statistics
9      means.append(bn_layer.running_mean.clone())
10     vars.append(bn_layer.running_var.clone())
11     model(x)
12
13 # Plot variance of batch means over time
14 plt.plot([m.std().item() for m in means])
15 plt.title('Variance of Batch Means')
16 plt.show()

```

9.6 Key Takeaways

Why normalization helps:

- Reduces internal covariate shift (distribution changes between layers)
- Smooths optimization landscape (allows higher learning rates)
- Provides regularization (slight noise from batch statistics)
- Reduces sensitivity to initialization

Batch Normalization:

- Normalizes over batch dimension
- Requires large enough batches (≥ 8)
- Different behavior in train/eval modes (critical!)
- Best for CNNs and MLPs with large batches
- Can speed up training by 2-10×

Layer Normalization:

- Normalizes over feature dimension
- Works with any batch size (even 1)
- Same behavior in train/eval modes
- Best for RNNs, Transformers, sequences
- Standard in modern NLP models

Implementation best practices:

- Place after Conv/Linear, before activation (standard)
- Set `bias=False` in layer before batch norm
- Always call `model.eval()` during evaluation
- Use `torch.no_grad()` during inference
- For sequences: prefer Layer Norm
- For small batches: use Layer Norm or Group Norm

Common mistakes:

- Forgetting `model.eval()` (very common bug!)
- Using Batch Norm with tiny batches
- Not setting `bias=False` (wastes parameters)
- Wrong normalization type for task (Batch Norm for RNNs)
- Not understanding train vs eval mode difference

Choosing normalization:

- **Default for CNNs:** Batch Norm
- **Default for Transformers/RNNs:** Layer Norm
- **Small batches (<8):** Layer Norm or Group Norm
- **Style transfer:** Instance Norm
- **Object detection:** Group Norm (common in modern detectors)

When NOT to use normalization:

- Very shallow networks (<5 layers)
- Already training stably

- Certain GANs (can cause mode collapse)
- When exact reproducibility needed (batch stats introduce randomness)

10 Recurrent Neural Networks (RNNs)

10.1 Introduction: Processing Sequential Data

So far, we've seen networks that process fixed-size inputs. But many real-world problems involve **sequences**:

- Time series (stock prices, sensor data, climate data)
- Text (words, characters, sentences)
- Audio (speech, music)
- Video (frames over time)
- Biological sequences (DNA, proteins)

Challenges with sequences:

1. Variable length (sentences have different lengths)
2. Temporal dependencies (past affects future)
3. Need to maintain state (remember what happened before)

Why not use MLPs or CNNs?

MLPs:

- Fixed input size (can't handle variable-length sequences)
- No notion of order (position 1 vs position 100 treated same)
- Can't share parameters across time (learn same pattern at different positions)

CNNs:

- Can handle sequences with 1D convolutions
- Limited temporal context (receptive field)
- No true memory of long-term dependencies

RNNs solve these problems by maintaining a **hidden state** that gets updated at each time step.

10.2 Theory: How RNNs Work

10.2.1 The Basic RNN

An RNN processes a sequence one element at a time, maintaining a hidden state:

At each time step t :

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

$$y_t = W_{hy}h_t + b_y$$

where:

- x_t : Input at time t
- h_t : Hidden state at time t (memory)
- y_t : Output at time t
- W_{hh} : Hidden-to-hidden weights (memory transformation)
- W_{xh} : Input-to-hidden weights
- W_{hy} : Hidden-to-output weights

Key idea: The hidden state h_t summarizes information from all previous time steps.

[The Hidden State] The hidden state h_t is the RNN's **memory**:

- h_0 : Initial memory (usually zeros)
- $h_1 = f(x_1, h_0)$: Combines first input with initial state
- $h_2 = f(x_2, h_1)$: Combines second input with previous state
- $h_3 = f(x_3, h_2)$: And so on...

By the end: h_T contains information about the entire sequence x_1, x_2, \dots, x_T .

Unrolled view:

$x_1 \rightarrow [\text{RNN}] \rightarrow h_1 \rightarrow [\text{RNN}] \rightarrow h_2 \rightarrow [\text{RNN}] \rightarrow h_3 \rightarrow \dots$
 \downarrow \downarrow \downarrow
 y_1 y_2 y_3

All boxes share the **same weights** (parameter sharing across time).

10.2.2 Parameter Sharing

Why parameter sharing matters:

For a sequence of length 100:

- MLP: Would need 100 separate weight matrices (millions of parameters)
- RNN: Uses same weights at each step (thousands of parameters)

The network learns patterns that work at **any position** in the sequence.

10.2.3 Sequence Modeling Patterns

1. One-to-Many (sequence generation):

$$\begin{array}{ccccccc} x & \rightarrow & [\text{RNN}] & \rightarrow & [\text{RNN}] & \rightarrow & [\text{RNN}] & \rightarrow & \dots \\ & & \downarrow & & \downarrow & & \downarrow & & \\ & & y_1 & & y_2 & & y_3 & & \end{array}$$

Example: Image captioning (image \rightarrow sentence)

2. Many-to-One (sequence classification):

$$\begin{array}{ccccccc} x_1 & \rightarrow & [\text{RNN}] & \rightarrow & x_2 & \rightarrow & [\text{RNN}] & \rightarrow & x_3 & \rightarrow & [\text{RNN}] \\ & & & & & & \downarrow & & & & \\ & & & & & & y & & & & \end{array}$$

Example: Sentiment analysis (sentence \rightarrow positive/negative)

3. Many-to-Many (same length):

$$\begin{array}{ccccccc} x_1 & \rightarrow & [\text{RNN}] & \rightarrow & x_2 & \rightarrow & [\text{RNN}] & \rightarrow & x_3 & \rightarrow & [\text{RNN}] \\ & & \downarrow & & \downarrow & & \downarrow & & & & \\ & & y_1 & & y_2 & & y_3 & & & & \end{array}$$

Example: Part-of-speech tagging (word \rightarrow tag)

4. Many-to-Many (different length, encoder-decoder):

$$\begin{array}{ccccccc} \text{Encoder: } & x_1 & \rightarrow & [\text{RNN}] & \rightarrow & x_2 & \rightarrow & [\text{RNN}] & \rightarrow & h \\ \text{Decoder: } & h & \rightarrow & [\text{RNN}] & \rightarrow & [\text{RNN}] & \rightarrow & [\text{RNN}] & & \\ & & & \downarrow & & \downarrow & & \downarrow & & \\ & & & y_1 & & y_2 & & y_3 & & \end{array}$$

Example: Machine translation (English \rightarrow French)

10.2.4 The Vanishing Gradient Problem

Problem: Vanilla RNNs struggle with long sequences.

When backpropagating through time, gradients are multiplied by W_{hh} at each step:

$$\frac{\partial h_t}{\partial h_0} = \prod_{i=1}^t W_{hh} \cdot \text{diag}(\tanh'(h_i))$$

If $|W_{hh}| < 1$: Gradients vanish exponentially (network can't learn long-term dependencies)

If $|W_{hh}| > 1$: Gradients explode (training becomes unstable)

In practice: Vanilla RNNs can only remember 10-20 steps back.

Solution: LSTM and GRU architectures (next sections).

10.3 Theory: Long Short-Term Memory (LSTM)

LSTMs solve the vanishing gradient problem through a clever gating mechanism.

10.3.1 LSTM Architecture

LSTM has two states:

- h_t : Hidden state (short-term memory)
- c_t : Cell state (long-term memory)

And three gates:

- f_t : Forget gate (what to remove from memory)
- i_t : Input gate (what new information to add)
- o_t : Output gate (what to output)

LSTM Equations:

1. Forget gate (what to forget):

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$$

2. Input gate (what to add):

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_c[h_{t-1}, x_t] + b_c) \quad (\text{candidate values})$$

3. Update cell state:

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

4. Output gate (what to output):

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \odot \tanh(c_t)$$

where \odot is element-wise multiplication, σ is sigmoid, $[a, b]$ is concatenation.

[LSTM Intuition] **Cell state (c_t):** The "memory highway"

- Runs through the entire sequence
- Modified only by additions and element-wise multiplications
- Gradients flow easily (no repeated matrix multiplications!)

Forget gate (f_t): "Should I forget the past?"

- Values near 0: Forget most of c_{t-1}
- Values near 1: Keep most of c_{t-1}

Input gate (i_t): "Should I add this new information?"

- Controls how much of \tilde{c}_t to add to memory

Output gate (o_t): "How much of memory should I expose?"

- Controls how much of c_t becomes h_t

Why LSTM works:

The cell state c_t provides a path where gradients can flow without repeated matrix multiplications (just element-wise operations). This solves vanishing gradients!

10.4 Theory: Gated Recurrent Unit (GRU)

GRU is a simplified version of LSTM with fewer parameters.

10.4.1 GRU Architecture

GRU combines forget and input gates into a single **update gate** and merges cell state with hidden state.

GRU Equations:

1. **Update gate (how much to update):**

$$z_t = \sigma(W_z[h_{t-1}, x_t] + b_z)$$

2. **Reset gate (how much past to forget):**

$$r_t = \sigma(W_r[h_{t-1}, x_t] + b_r)$$

3. **Candidate hidden state:**

$$\tilde{h}_t = \tanh(W_h[r_t \odot h_{t-1}, x_t] + b_h)$$

4. **Update hidden state:**

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

Intuition:

- $z_t \approx 0$: Keep old state h_{t-1} (don't update)
- $z_t \approx 1$: Use new candidate \tilde{h}_t (full update)
- $z_t \approx 0.5$: Blend old and new

LSTM vs GRU:

Aspect	LSTM	GRU
Gates	3 (forget, input, output)	2 (update, reset)
States	2 (h_t , c_t)	1 (h_t)
Parameters	More	Fewer (faster)
Performance	Slightly better	Nearly as good
Training speed	Slower	Faster
Memory	More	Less

Table 8: LSTM vs GRU comparison

When to use each:

- **LSTM:** Default choice, especially for complex tasks
- **GRU:** When speed matters or dataset is smaller
- **Rule of thumb:** Try both, see which works better

10.5 Implementation: RNNs in PyTorch

10.5.1 Vanilla RNN

```

1 import torch
2 import torch.nn as nn
3
4 # Simple RNN cell
5 rnn_cell = nn.RNNCell(input_size=10, hidden_size=20)
6
7 # Input: (batch, input_size)
8 # Hidden: (batch, hidden_size)
9 x = torch.randn(32, 10) # Batch of 32
10 h = torch.zeros(32, 20) # Initial hidden state
11
12 # One step
13 h_next = rnn_cell(x, h)
14 print(h_next.shape) # torch.Size([32, 20])
15
16 # Full RNN layer (processes entire sequence)
17 rnn = nn.RNN(input_size=10, hidden_size=20, num_layers=1,
18               batch_first=True)
19
20 # Input: (batch, seq_len, input_size)
21 x = torch.randn(32, 100, 10) # 32 sequences, length 100
22 h0 = torch.zeros(1, 32, 20) # (num_layers, batch,
23                               hidden_size)
24
25 # Forward pass
26 output, h_n = rnn(x, h0)
27 # output: (32, 100, 20) - outputs at each time step
28 # h_n: (1, 32, 20) - final hidden state
29 print(output.shape, h_n.shape)

```

[batch_first Parameter] PyTorch RNNs default to shape (seq_len, batch, features), but this is confusing!

Always use `batch_first=True` for shape (batch, seq_len, features):

```

1 rnn = nn.RNN(..., batch_first=True) # Recommended
2 # Input: (batch, seq_len, features)
3 # Output: (batch, seq_len, hidden_size)

```

10.5.2 LSTM

```

1  # LSTM cell
2  lstm_cell = nn.LSTMCell(input_size=10, hidden_size=20)
3
4  x = torch.randn(32, 10)
5  h = torch.zeros(32, 20)
6  c = torch.zeros(32, 20)  # Cell state
7
8  # One step
9  h_next, c_next = lstm_cell(x, (h, c))  # Note: tuple of (h, c)
10
11 # Full LSTM layer
12 lstm = nn.LSTM(input_size=10, hidden_size=20, num_layers=2,
13                batch_first=True, dropout=0.2)
14
15 x = torch.randn(32, 100, 10)
16 h0 = torch.zeros(2, 32, 20)  # (num_layers, batch, hidden)
17 c0 = torch.zeros(2, 32, 20)  # (num_layers, batch, hidden)
18
19 # Forward pass
20 output, (h_n, c_n) = lstm(x, (h0, c0))
21 # output: (32, 100, 20) - outputs at each time step
22 # h_n: (2, 32, 20) - final hidden state for each layer
23 # c_n: (2, 32, 20) - final cell state for each layer
24
25 print(output.shape, h_n.shape, c_n.shape)

```

[LSTM Hidden State is a Tuple!] LSTM returns (h, c) as a tuple, not just h:

```

1  # WRONG:
2  output, hidden = lstm(x)
3  h_n = hidden  # This is a tuple, not a tensor!
4
5  # CORRECT:
6  output, (h_n, c_n) = lstm(x, (h0, c0))
7  # or
8  output, hidden = lstm(x, (h0, c0))
9  h_n, c_n = hidden

```

10.5.3 GRU

```

1  # GRU (similar to LSTM, but only one hidden state)
2  gru = nn.GRU(input_size=10, hidden_size=20, num_layers=2,
3               batch_first=True, dropout=0.2)
4
5  x = torch.randn(32, 100, 10)
6  h0 = torch.zeros(2, 32, 20)
7
8  # Forward pass
9  output, h_n = gru(x, h0)  # Only h, no cell state
10 # output: (32, 100, 20)
11 # h_n: (2, 32, 20)
12
13 print(output.shape, h_n.shape)

```

10.5.4 Building a Sequence Classifier

```

1 class SequenceClassifier(nn.Module):
2     """Many-to-one: sequence      single label."""
3
4     def __init__(self, input_size, hidden_size, num_classes,
5 num_layers=2):
6         super().__init__()
7
8         self.lstm = nn.LSTM(
9             input_size=input_size,
10            hidden_size=hidden_size,
11            num_layers=num_layers,
12            batch_first=True,
13            dropout=0.2 if num_layers > 1 else 0
14        )
15
16        self.fc = nn.Linear(hidden_size, num_classes)
17
18    def forward(self, x):
19        # x: (batch, seq_len, input_size)
20
21        # LSTM forward
22        output, (h_n, c_n) = self.lstm(x)
23        # output: (batch, seq_len, hidden_size)
24        # h_n: (num_layers, batch, hidden_size)
25
26        # Use last time step's output
27        last_output = output[:, -1, :] # (batch, hidden_size)
28
29        # Or use final hidden state from last layer
30        # last_hidden = h_n[-1] # (batch, hidden_size)
31
32        # Classify
33        logits = self.fc(last_output)
34        return logits
35
36 # Test
37 model = SequenceClassifier(input_size=10, hidden_size=64,
38 num_classes=3)
39 x = torch.randn(32, 100, 10) # 32 sequences, length 100
40 output = model(x)
41 print(output.shape) # torch.Size([32, 3])

```


10.5.5 Sequence-to-Sequence (Many-to-Many)

```

1 class SequenceLabeler(nn.Module):
2     """Many-to-many: each input      each output (same length)
3     ."""
4
5     def __init__(self, input_size, hidden_size, num_classes):
6         super().__init__()
7
8         self.lstm = nn.LSTM(
9             input_size=input_size,
10            hidden_size=hidden_size,
11            num_layers=2,
12            batch_first=True,
13            dropout=0.2
14        )
15
16        # Apply classifier at each time step
17        self.fc = nn.Linear(hidden_size, num_classes)
18
19    def forward(self, x):
20        # x: (batch, seq_len, input_size)
21
22        output, _ = self.lstm(x)
23        # output: (batch, seq_len, hidden_size)
24
25        # Apply classifier to all time steps
26        logits = self.fc(output)
27        # logits: (batch, seq_len, num_classes)
28
29        return logits
30
31    # Test
32    model = SequenceLabeler(input_size=10, hidden_size=64,
33                            num_classes=5)
34    x = torch.randn(32, 100, 10)
35    output = model(x)
36    print(output.shape) # torch.Size([32, 100, 5])

```

10.5.6 Bidirectional RNNs

Process sequence in both directions (forward and backward):

```

1 class BidirectionalRNN(nn.Module):
2     """Bidirectional LSTM for sequence classification."""
3
4     def __init__(self, input_size, hidden_size, num_classes):
5         super().__init__()
6
7         self.lstm = nn.LSTM(
8             input_size=input_size,
9             hidden_size=hidden_size,
10            num_layers=2,
11            batch_first=True,
12            dropout=0.2,
13            bidirectional=True # Process both directions
14        )
15
16        # Hidden size is doubled (forward + backward)
17        self.fc = nn.Linear(hidden_size * 2, num_classes)
18
19    def forward(self, x):

```

```
20     # x: (batch, seq_len, input_size)
21
22     output, _ = self.lstm(x)
23     # output: (batch, seq_len, hidden_size * 2)
24
25     # Use last time step (or all for seq-to-seq)
26     last_output = output[:, -1, :]
27
28     logits = self.fc(last_output)
29     return logits
30
31 # Test
32 model = BidirectionalRNN(input_size=10, hidden_size=64,
33                           num_classes=3)
33 x = torch.randn(32, 100, 10)
34 output = model(x)
35 print(output.shape) # torch.Size([32, 3])
```

When to use bidirectional:

- Text classification (can see future context)
- Part-of-speech tagging
- Named entity recognition
- Any task where you have the full sequence before prediction

Don't use bidirectional for:

- Time series prediction (can't see future!)
- Language generation (need causal/autoregressive)
- Online/streaming applications

10.5.7 Handling Variable-Length Sequences

```

1 from torch.nn.utils.rnn import pack_padded_sequence,
  pad_packed_sequence
2
3 class VariableLengthRNN(nn.Module):
4     """Efficiently handle sequences of different lengths."""
5
6     def __init__(self, input_size, hidden_size, num_classes):
7         super().__init__()
8         self.lstm = nn.LSTM(input_size, hidden_size,
9                               batch_first=True)
10        self.fc = nn.Linear(hidden_size, num_classes)
11
12    def forward(self, x, lengths):
13        """
14        Args:
15        x: Padded sequences (batch, max_len, input_size)
16        lengths: Actual lengths of each sequence
17        """
18        # Pack sequences (skip padding in computation)
19        packed = pack_padded_sequence(x, lengths, batch_first
20                                     =True,
21                                     enforce_sorted=False)
22
23        # LSTM forward on packed sequence
24        packed_output, (h_n, c_n) = self.lstm(packed)
25
26        # Unpack
27        output, _ = pad_packed_sequence(packed_output,
28                                       batch_first=True)
29
30        # Use final hidden state
31        logits = self.fc(h_n[-1])
32
33        return logits
34
35    # Example usage
36    # Sequences of different lengths: [50, 30, 40]
37    lengths = torch.tensor([50, 30, 40])
38    max_len = lengths.max()
39
40    # Pad to max length
41    x = torch.randn(3, max_len, 10) # Some values are padding
42
43    model = VariableLengthRNN(10, 64, 3)
44    output = model(x, lengths)
45    print(output.shape) # torch.Size([3, 3])

```

Why pack sequences?

- Skip computation on padding (faster)
- More accurate (padding doesn't affect hidden state)
- More memory efficient

10.5.8 Gradient Clipping (Essential for RNNs)

RNNs are prone to exploding gradients. Gradient clipping is **essential**.

```

1  # Method 1: Clip by norm (recommended)
2  torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm
   =5.0)
3
4  # Method 2: Clip by value
5  torch.nn.utils.clip_grad_value_(model.parameters(),
   clip_value=0.5)
6
7  # In training loop
8  for epoch in range(num_epochs):
9      for batch in dataloader:
10         optimizer.zero_grad()
11
12         output = model(batch)
13         loss = criterion(output, target)
14         loss.backward()
15
16         # Clip gradients AFTER backward, BEFORE step
17         torch.nn.utils.clip_grad_norm_(model.parameters(),
18         max_norm=5.0)
19
20         optimizer.step()

```

[Always Use Gradient Clipping for RNNs] Without gradient clipping:

- Loss may suddenly spike to NaN
- Training becomes unstable
- Network may never converge

Typical values:

- max_norm=1.0: Conservative (very stable)
- max_norm=5.0: Standard (good default)
- max_norm=10.0: Aggressive (faster but less stable)

10.5.9 Training Tips for RNNs

1. Learning rate

```

1  # RNNs typically need lower learning rates than CNNs
2  optimizer = torch.optim.Adam(model.parameters(), lr=0.001)  #
   Good start
3
4  # Use learning rate scheduling
5  scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(
6      optimizer, mode='min', factor=0.5, patience=5
7  )

```

2. Weight initialization

```

1  def init_lstm_weights(lstm):
2      """Initialize LSTM weights for better convergence."""
3      for name, param in lstm.named_parameters():
4          if 'weight_ih' in name:
5              # Input-hidden weights: Xavier
6              nn.init.xavier_uniform_(param)

```

```

7         elif 'weight_hh' in name:
8             # Hidden-hidden weights: Orthogonal (better for
              RNNs)
9             nn.init.orthogonal_(param)
10        elif 'bias' in name:
11            # Biases: zeros, except forget gate bias = 1
12            nn.init.zeros_(param)
13            # Set forget gate bias to 1 (helps gradient flow)
14            n = param.size(0)
15            param.data[n//4:n//2].fill_(1.0)
16
17    # Apply to model
18    for name, module in model.named_modules():
19        if isinstance(module, nn.LSTM):
20            init_lstm_weights(module)

```

3. Dropout placement

```

1  # Dropout between layers (built-in)
2  lstm = nn.LSTM(input_size, hidden_size, num_layers=3, dropout
              =0.2)
3  # This applies dropout between layers, not within recurrent
              connections
4
5  # For dropout on inputs/outputs:
6  class RNNWithDropout(nn.Module):
7      def __init__(self, input_size, hidden_size, num_classes):
8          super().__init__()
9          self.input_dropout = nn.Dropout(0.2)
10         self.lstm = nn.LSTM(input_size, hidden_size,
              num_layers=2)
11         self.output_dropout = nn.Dropout(0.5)
12         self.fc = nn.Linear(hidden_size, num_classes)
13
14     def forward(self, x):
15         x = self.input_dropout(x)
16         output, _ = self.lstm(x)
17         output = self.output_dropout(output[:, -1, :])
18         return self.fc(output)

```

4. Batch size considerations

RNNs benefit from larger batch sizes:

- Batch size 32-128 typically good
- Too small (<16): Noisy gradients, unstable training
- Too large (>256): May hurt generalization

10.5.10 Common RNN Architectures

Stacked RNNs (Deep RNNs):

```

1 class DeepLSTM(nn.Module):
2     """Multiple LSTM layers stacked."""
3
4     def __init__(self, input_size, hidden_size, num_layers,
5 num_classes):
6         super().__init__()
7
8         # num_layers > 1 creates stacked LSTM
9         self.lstm = nn.LSTM(
10             input_size=input_size,
11             hidden_size=hidden_size,
12             num_layers=num_layers,
13             batch_first=True,
14             dropout=0.3 if num_layers > 1 else 0 # Dropout
15             between layers
16         )
17
18         self.fc = nn.Linear(hidden_size, num_classes)
19
20     def forward(self, x):
21         output, (h_n, c_n) = self.lstm(x)
22         # h_n[-1]: hidden state from the last layer
23         return self.fc(h_n[-1])
24
25 # Typical depths: 2-4 layers
26 # Deeper than 4 often doesn't help much
27 model = DeepLSTM(input_size=10, hidden_size=128, num_layers
28 =3,
29 num_classes=10)

```

Encoder-Decoder (Seq2Seq):

```

1 class Seq2Seq(nn.Module):
2     """Encoder-decoder for sequence-to-sequence tasks."""
3
4     def __init__(self, input_size, hidden_size, output_size):
5         super().__init__()
6
7         # Encoder: processes input sequence
8         self.encoder = nn.LSTM(input_size, hidden_size,
9 batch_first=True)
10
11         # Decoder: generates output sequence
12         self.decoder = nn.LSTM(output_size, hidden_size,
13 batch_first=True)
14
15         # Output projection
16         self.fc = nn.Linear(hidden_size, output_size)
17
18     def forward(self, src, tgt, teacher_forcing_ratio=0.5):
19         """
20         Args:
21         src: Source sequence (batch, src_len, input_size)
22         tgt: Target sequence (batch, tgt_len, output_size)
23         teacher_forcing_ratio: Probability of using true
24         target vs prediction
25         """

```

```

23     batch_size = src.size(0)
24     tgt_len = tgt.size(1)
25     output_size = tgt.size(2)
26
27     # Encode source
28     _, (h_n, c_n) = self.encoder(src)
29
30     # Initialize decoder hidden state with encoder final
state
31     decoder_hidden = (h_n, c_n)
32
33     # First decoder input: start token (zeros)
34     decoder_input = torch.zeros(batch_size, 1,
35                                output_size).to(src.device)
36
37     outputs = []
38
39     # Generate sequence step by step
40     for t in range(tgt_len):
41         # Decode one step
42         decoder_output, decoder_hidden = self.decoder(
43             decoder_input, decoder_hidden
44         )
45
46         # Project to output space
47         prediction = self.fc(decoder_output)
48         outputs.append(prediction)
49
50         # Teacher forcing: use true target or prediction?
51         use_teacher_forcing = torch.rand(1).item() <
52         teacher_forcing_ratio
53
54         if use_teacher_forcing:
55             decoder_input = tgt[:, t:t+1, :] # Use true
56             target
57         else:
58             decoder_input = prediction # Use prediction
59
60         # Stack all outputs
61         outputs = torch.cat(outputs, dim=1)
62     return outputs
63
64 # Test
65 model = Seq2Seq(input_size=10, hidden_size=64, output_size
66                =10)
67 src = torch.randn(32, 20, 10) # 32 sequences, length 20
68 tgt = torch.randn(32, 15, 10) # Target length 15
69 output = model(src, tgt)
70 print(output.shape) # torch.Size([32, 15, 10])

```

10.5.11 Time Series Forecasting Example

```

1 class TimeSeriesForecaster(nn.Module):
2     """Predict future values from past observations."""
3
4     def __init__(self, input_size=1, hidden_size=64,
5                 num_layers=2, forecast_horizon=10):
6         super().__init__()
7
8         self.hidden_size = hidden_size
9         self.num_layers = num_layers
10        self.forecast_horizon = forecast_horizon
11
12        # Encoder LSTM
13        self.lstm = nn.LSTM(
14            input_size=input_size,
15            hidden_size=hidden_size,
16            num_layers=num_layers,
17            batch_first=True,
18            dropout=0.2 if num_layers > 1 else 0
19        )
20
21        # Decoder: predict next step
22        self.fc = nn.Linear(hidden_size, input_size)
23
24    def forward(self, x, future_steps=0):
25        """
26        Args:
27            x: Past observations (batch, seq_len, input_size)
28            future_steps: How many steps to predict into
29            future
26
27        Returns:
28            predictions: (batch, seq_len + future_steps,
29            input_size)
30        """
31        batch_size = x.size(0)
32
33        # Process historical data
34        output, (h, c) = self.lstm(x)
35
36        # Predict on historical data
37        predictions = self.fc(output)
38
39        # If we need to predict future steps
40        if future_steps > 0:
41            future_preds = []
42
43            # Use last prediction as next input
44            last_pred = predictions[:, -1:, :]
45
46            for _ in range(future_steps):
47                # Predict next step
48                output, (h, c) = self.lstm(last_pred, (h, c))
49                pred = self.fc(output)
50                future_preds.append(pred)
51
52            # Use prediction as next input (
53            autoregressive)
54            last_pred = pred
55
56

```



```

57         # Concatenate all predictions
58         future_preds = torch.cat(future_preds, dim=1)
59         predictions = torch.cat([predictions,
        future_preds], dim=1)
60
61         return predictions
62
63     # Example usage
64     model = TimeSeriesForecaster(input_size=1, hidden_size=64)
65
66     # Historical data: 100 time steps
67     x = torch.randn(32, 100, 1)
68
69     # Predict on historical + 20 future steps
70     predictions = model(x, future_steps=20)
71     print(predictions.shape) # torch.Size([32, 120, 1])

```

10.6 Exercises

10.1: Simple Sequence Classification - ★★

Goal: Build your first RNN.

1. Generate synthetic sequences: sine waves (class 0) and cosine waves (class 1)
2. Each sequence: 50 time steps
3. Build an LSTM classifier
4. Train for 20 epochs
5. Achieve >95% accuracy

Starter code:

```

1  import numpy as np
2
3  def generate_sine_data(n_samples=1000, seq_len=50):
4      """Generate sine wave sequences."""
5      X = []
6      y = []
7      for _ in range(n_samples):
8          # Random frequency and phase
9          freq = np.random.uniform(0.5, 2.0)
10         phase = np.random.uniform(0, 2*np.pi)
11         t = np.linspace(0, 4*np.pi, seq_len)
12
13         if np.random.rand() < 0.5:
14             # Sine wave
15             seq = np.sin(freq * t + phase)
16             label = 0
17         else:
18             # Cosine wave
19             seq = np.cos(freq * t + phase)
20             label = 1
21
22         X.append(seq)
23         y.append(label)
24
25     return torch.FloatTensor(X).unsqueeze(-1), torch.
        LongTensor(y)
26
27     # Your LSTM classifier here

```

10.2: LSTM vs GRU Comparison - ★★★**Goal:** Compare LSTM and GRU empirically.

1. Generate a moderately complex sequence task
2. Train identical architectures with:
 - LSTM (2 layers, 128 hidden)
 - GRU (2 layers, 128 hidden)
 - Vanilla RNN (2 layers, 128 hidden)
3. Compare:
 - Training speed (time per epoch)
 - Convergence speed (epochs to 90% accuracy)
 - Final accuracy
 - Number of parameters

Expected results:

- LSTM and GRU should perform similarly
- GRU should be faster (fewer parameters)
- Vanilla RNN should struggle (vanishing gradients)

10.3: Sequence-to-Sequence Labeling - ★★★**Goal:** Many-to-many prediction.

1. Task: Given a noisy sine wave, denoise it at each time step
2. Input: sine + Gaussian noise
3. Target: clean sine wave
4. Build a bidirectional LSTM
5. Output prediction at each time step
6. Visualize: noisy input vs clean output

Evaluation:

- Mean squared error on test set
- Visual inspection of denoised signals

10.4: Variable-Length Sequences - ★★★**Goal:** Handle sequences of different lengths efficiently.

1. Generate sequences with random lengths (20-100)
2. Implement proper padding and packing:
 - Pad sequences to same length
 - Use `pack_padded_sequence`
 - Use `pad_packed_sequence`
3. Train two models:
 - Without packing (processes padding)
 - With packing (skips padding)
4. Compare:
 - Training time
 - Accuracy

- Verify packing gives same results

10.5: Time Series Forecasting - ★★★★★

Goal: Predict future values.

1. Use real time series data (or generate complex synthetic data)
2. Sliding window approach:
 - Input: past 50 steps
 - Target: next 10 steps
3. Build forecasting model with LSTM
4. Implement autoregressive prediction (use predictions as inputs)
5. Evaluate:
 - One-step-ahead prediction
 - Multi-step-ahead prediction
 - Compare with baseline (simple moving average)
6. Visualize predictions vs ground truth

Metrics:

```

1 def calculate_metrics(y_true, y_pred):
2     """Calculate forecasting metrics."""
3     mse = torch.mean((y_true - y_pred) ** 2)
4     mae = torch.mean(torch.abs(y_true - y_pred))
5
6     # MAPE (Mean Absolute Percentage Error)
7     mape = torch.mean(torch.abs((y_true - y_pred) / (
8         y_true + 1e-8))) * 100
9
10    return {
11        'MSE': mse.item(),
12        'MAE': mae.item(),
13        'MAPE': mape.item()
14    }
```

10.6: Gradient Clipping Investigation - ★★★★★

Goal: Understand why gradient clipping is essential.

1. Train an RNN on a long sequence task (100+ steps)
2. Train multiple times with different settings:
 - No gradient clipping
 - max_norm=1.0
 - max_norm=5.0
 - max_norm=10.0
3. Monitor and plot:
 - Gradient norms over time
 - Loss curves
 - Percentage of batches where gradients are clipped
4. Observe:
 - Without clipping: loss spikes, NaN values
 - With clipping: stable training

Visualization code:

```

1 def track_gradient_norms(model):
2     """Track gradient norms during training."""
3     total_norm = 0
4     for p in model.parameters():
5         if p.grad is not None:
6             param_norm = p.grad.data.norm(2)
7             total_norm += param_norm.item() ** 2
8     total_norm = total_norm ** 0.5
9     return total_norm
10
11 # In training loop
12 gradient_norms = []
13 for epoch in range(num_epochs):
14     for batch in dataloader:
15         optimizer.zero_grad()
16         loss = ...
17         loss.backward()
18
19         # Track before clipping
20         norm_before = track_gradient_norms(model)
21         gradient_norms.append(norm_before)
22
23         torch.nn.utils.clip_grad_norm_(model.parameters
24         (), max_norm=5.0)
25         optimizer.step()
26
27 # Plot
28 plt.plot(gradient_norms)
29 plt.axhline(y=5.0, color='r', linestyle='--', label='
30             Clip threshold')
31 plt.xlabel('Step')
32 plt.ylabel('Gradient Norm')
33 plt.legend()
34 plt.show()

```

10.7 Key Takeaways

When to use RNNs:

- Sequential data with temporal dependencies
- Variable-length sequences
- Time series forecasting
- Text processing (though Transformers now dominate)
- When order matters

Architecture choices:

- **Vanilla RNN:** Almost never use (vanishing gradients)
- **LSTM:** Default choice, especially for complex tasks
- **GRU:** When speed matters or dataset is smaller
- **Bidirectional:** When you have full sequence (not streaming)
- **Stacked (2-4 layers):** For complex patterns

Essential training practices:

- **Always use gradient clipping** (`max_norm=5.0` is good default)
- Initialize forget gate bias to 1 (helps gradient flow)
- Use orthogonal initialization for hidden-to-hidden weights
- Lower learning rates than CNNs (0.001 typical)
- Larger batch sizes (32-128)
- Pack sequences when lengths vary (efficiency)

Common patterns:

- **Many-to-one:** Classification (use last hidden state)
- **Many-to-many (same length):** Sequence labeling (output at each step)
- **Many-to-many (different length):** Seq2seq with encoder-decoder
- **One-to-many:** Generation (feed output back as input)

Why LSTMs work better than vanilla RNNs:

- Cell state provides gradient highway (no repeated matrix multiplications)
- Gates control information flow (selective memory)
- Can learn dependencies over 100+ steps
- Forget gate prevents exploding activations

Limitations of RNNs:

- Sequential processing (slow, can't parallelize)
- Still struggle with very long sequences (>1000 steps)
- Transformers now preferred for many NLP tasks
- Require careful tuning (gradient clipping, learning rate)

Practical tips:

- Start with 2-layer LSTM, hidden size = 128 or 256
- Use `batch_first=True` always
- Remember LSTM returns `(h, c)` tuple
- Dropout between layers (0.2-0.3), higher on outputs (0.5)

- Monitor gradient norms (should be <10)
- Use early stopping (RNNs prone to overfitting)

Debugging checklist:

- Loss explodes? → Add/increase gradient clipping
- Loss stuck? → Check learning rate, try lower
- Poor performance? → Try LSTM instead of GRU, add layers
- Slow training? → Use packing for variable lengths
- Different train/test? → Check dropout mode

Modern alternatives:

- **Transformers:** Better for most NLP tasks (parallelizable)
- **Temporal CNNs:** Good for some time series (efficient)
- **State space models:** Emerging alternative (S4, Mamba)

However, RNNs still useful for:

- Online/streaming applications (Transformers need full context)
- Very long sequences where Transformers too expensive
- Small datasets (fewer parameters than Transformers)
- Scientific time series (physics-informed RNNs)

11 Attention & Transformers

11.1 Introduction: Beyond RNNs

RNNs revolutionized sequence modeling, but they have fundamental limitations:

Problems with RNNs:

1. **Sequential processing:** Can't parallelize (slow!)
2. **Long-range dependencies:** Even LSTM struggles beyond 100 steps
3. **Information bottleneck:** All info must pass through hidden state
4. **Vanishing gradients:** Still an issue for very long sequences

The **Attention mechanism** (2014) and **Transformers** (2017) solved these problems.

Key innovation: Instead of compressing everything into a fixed-size hidden state, let the model **attend to** any part of the input directly.

Impact:

- GPT, BERT, ChatGPT: All based on Transformers
- Transformers now dominate NLP, vision (ViT), multimodal (CLIP), and more
- Enabled training on massive datasets (parallelization)
- State-of-the-art on virtually all sequence tasks

11.2 Theory: The Attention Mechanism

11.2.1 Intuition: What is Attention?

Imagine reading a document and answering a question. You don't re-read the entire document—you **focus** on relevant parts.

Example:

- Document: "The cat sat on the mat. The dog played in the garden."
- Query: "Where is the cat?"
- You attend to: "The cat sat on the mat" (relevant)
- You ignore: "The dog played in the garden" (irrelevant)

Attention does this automatically:

- Given a query, compute relevance scores for all inputs
- Higher scores = more attention
- Output is weighted sum of values, weighted by attention scores

11.2.2 Mathematical Formulation

Attention has three components:

- **Query (Q):** What we're looking for
- **Keys (K):** What each position represents
- **Values (V):** The actual content at each position

Step 1: Compute attention scores

$$\text{scores} = QK^T$$

Each element (i, j) is how relevant key j is to query i .

Step 2: Normalize with softmax

$$\text{attention_weights} = \text{softmax}(\text{scores})$$

Converts scores to probabilities (sum to 1).

Step 3: Weighted sum of values

$$\text{output} = \text{attention_weights} \cdot V$$

Complete formula (scaled dot-product attention):

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

where d_k is the dimension of keys (scaling prevents large dot products).

[Why Scaling by $\sqrt{d_k}$?] Without scaling, dot products can become very large when d_k is large:

- Large dot products \rightarrow softmax saturates \rightarrow tiny gradients
- Scaling by $\sqrt{d_k}$ keeps values in reasonable range
- Ensures stable gradients

Example: If $d_k = 512$, dot products could be ~ 512 without scaling. After scaling: $\sim \sqrt{512} \approx 22.6$.

11.2.3 Self-Attention vs Cross-Attention

Self-Attention: Attention within same sequence

Query, Key, Value all come from the same source.

Use case: Understanding relationships between words in a sentence.

Example:

- Sentence: "The animal didn't cross the street because it was too tired"
- "it" attends to "animal" (not "street")

Cross-Attention: Attention between two sequences

Query from one sequence, Keys and Values from another.

Use case: Machine translation, encoder-decoder models.

Example:

- English (Keys/Values): "The cat is black"
- French decoder (Query): "Le chat est..."
- Decoder attends to relevant English words

11.2.4 Multi-Head Attention

Instead of one attention mechanism, use multiple in parallel:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where each head is:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

Why multiple heads?

- Different heads can focus on different aspects
- Head 1: Syntax (subject-verb agreement)
- Head 2: Semantics (meaning relationships)
- Head 3: Position (nearby words)
- Richer representations than single attention

Typical values:

- $h = 8$ or $h = 16$ heads
- $d_{\text{model}} = 512$ (total dimension)
- $d_k = d_v = d_{\text{model}}/h = 64$ per head

11.2.5 Positional Encoding

Problem: Attention has no notion of order!

The sentence "cat chased dog" and "dog chased cat" look identical to attention (same words, same attention scores).

Solution: Add positional information to input embeddings.

Sinusoidal positional encoding:

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$
$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

where pos is position, i is dimension index.

Properties:

- Different frequency for each dimension
- Can extrapolate to longer sequences than seen during training
- Model can learn to attend by relative position

Alternative: Learned positional embeddings (fixed maximum length).

11.3 Theory: The Transformer Architecture

11.3.1 Complete Architecture

Transformer consists of:

- **Encoder:** Processes input sequence
- **Decoder:** Generates output sequence

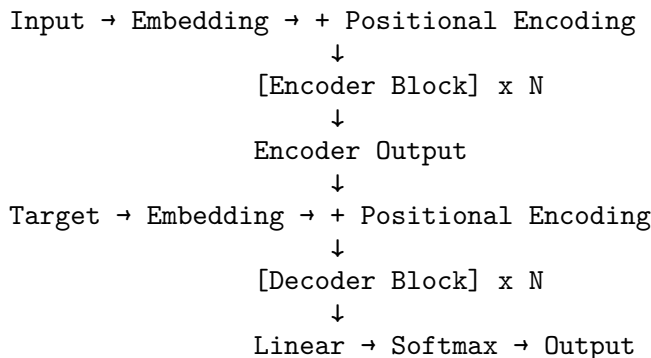
Encoder Block (repeated N times):

1. Multi-head self-attention
2. Add & Layer Norm (residual connection)
3. Feed-forward network (two linear layers with ReLU)
4. Add & Layer Norm (residual connection)

Decoder Block (repeated N times):

1. Masked multi-head self-attention (can't see future)
2. Add & Layer Norm
3. Multi-head cross-attention (attend to encoder output)
4. Add & Layer Norm
5. Feed-forward network
6. Add & Layer Norm

Diagram:



11.3.2 Why Transformers Work So Well

1. Parallelization

RNN: Must process step-by-step (slow)

Transformer: All positions processed in parallel (fast!)

Training on long sequences: 10-100× faster than RNNs.

2. Long-range dependencies

RNN: Information must flow through many steps (degrades)

Transformer: Direct connection between any two positions (one attention operation)

3. Scalability

Transformers scale beautifully with:

- More data (billions of tokens)

- More parameters (billions to trillions)
- More compute (hundreds of GPUs)

4. Transfer learning

Pre-train once on massive data → Fine-tune for specific tasks.

Examples: BERT, GPT, T5.

5. Inductive biases

Minimal assumptions about data structure:

- CNNs assume local structure (good for images)
- RNNs assume sequential processing (good for sequences)
- Transformers learn structure from data (flexible!)

11.4 Implementation: Building Attention from Scratch

11.4.1 Scaled Dot-Product Attention

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import math
5
6 def scaled_dot_product_attention(query, key, value, mask=None
7 ):
8     """
9     Scaled dot-product attention.
10
11     Args:
12         query: (batch, seq_len_q, d_k)
13         key: (batch, seq_len_k, d_k)
14         value: (batch, seq_len_k, d_v)
15         mask: Optional mask (batch, seq_len_q, seq_len_k)
16
17     Returns:
18         output: (batch, seq_len_q, d_v)
19         attention_weights: (batch, seq_len_q, seq_len_k)
20     """
21     d_k = query.size(-1)
22
23     # Compute attention scores: Q @ K^T
24     scores = torch.matmul(query, key.transpose(-2, -1)) # (
25     batch, seq_len_q, seq_len_k)
26
27     # Scale by sqrt(d_k)
28     scores = scores / math.sqrt(d_k)
29
30     # Apply mask if provided (set masked positions to -inf)
31     if mask is not None:
32         scores = scores.masked_fill(mask == 0, -1e9)
33
34     # Softmax to get attention weights
35     attention_weights = F.softmax(scores, dim=-1) # (batch,
36     seq_len_q, seq_len_k)
37
38     # Weighted sum of values
39     output = torch.matmul(attention_weights, value) # (batch
40     , seq_len_q, d_v)
41
42     return output, attention_weights
43
44 # Test
45 batch_size = 2
46 seq_len = 10
47 d_k = 64
48
49 Q = torch.randn(batch_size, seq_len, d_k)
50 K = torch.randn(batch_size, seq_len, d_k)
51 V = torch.randn(batch_size, seq_len, d_k)
52
53 output, weights = scaled_dot_product_attention(Q, K, V)
54 print(f"Output shape: {output.shape}") # (2, 10, 64)
55 print(f"Attention weights shape: {weights.shape}") # (2, 10,
56 10)
57 print(f"Weights sum to 1: {weights.sum(dim=-1)}") # All ones

```

11.4.2 Multi-Head Attention

```

1 class MultiHeadAttention(nn.Module):
2     """Multi-head attention mechanism."""
3
4     def __init__(self, d_model, num_heads, dropout=0.1):
5         """
6         Args:
7             d_model: Model dimension (e.g., 512)
8             num_heads: Number of attention heads (e.g., 8)
9             dropout: Dropout probability
10        """
11        super().__init__()
12
13        assert d_model % num_heads == 0, "d_model must be
divisible by num_heads"
14
15        self.d_model = d_model
16        self.num_heads = num_heads
17        self.d_k = d_model // num_heads # Dimension per head
18
19        # Linear projections for Q, K, V
20        self.W_q = nn.Linear(d_model, d_model)
21        self.W_k = nn.Linear(d_model, d_model)
22        self.W_v = nn.Linear(d_model, d_model)
23
24        # Output projection
25        self.W_o = nn.Linear(d_model, d_model)
26
27        self.dropout = nn.Dropout(dropout)
28
29    def split_heads(self, x):
30        """
31        Split last dimension into (num_heads, d_k).
32
33        Input: (batch, seq_len, d_model)
34        Output: (batch, num_heads, seq_len, d_k)
35        """
36        batch_size, seq_len, d_model = x.size()
37
38        # Reshape to (batch, seq_len, num_heads, d_k)
39        x = x.view(batch_size, seq_len, self.num_heads, self.
d_k)
40
41        # Transpose to (batch, num_heads, seq_len, d_k)
42        return x.transpose(1, 2)
43
44    def combine_heads(self, x):
45        """
46        Combine heads back.
47
48        Input: (batch, num_heads, seq_len, d_k)
49        Output: (batch, seq_len, d_model)
50        """
51        batch_size, num_heads, seq_len, d_k = x.size()
52
53        # Transpose to (batch, seq_len, num_heads, d_k)
54        x = x.transpose(1, 2)
55
56        # Reshape to (batch, seq_len, d_model)
57        return x.contiguous().view(batch_size, seq_len, self.
d_model)

```

```

58
59     def forward(self, query, key, value, mask=None):
60         """
61         Args:
62             query: (batch, seq_len_q, d_model)
63             key: (batch, seq_len_k, d_model)
64             value: (batch, seq_len_k, d_model)
65             mask: Optional (batch, seq_len_q, seq_len_k)
66         """
67         batch_size = query.size(0)
68
69         # Linear projections
70         Q = self.W_q(query) # (batch, seq_len_q, d_model)
71         K = self.W_k(key)   # (batch, seq_len_k, d_model)
72         V = self.W_v(value) # (batch, seq_len_k, d_model)
73
74         # Split into multiple heads
75         Q = self.split_heads(Q) # (batch, num_heads,
seq_len_q, d_k)
76         K = self.split_heads(K) # (batch, num_heads,
seq_len_k, d_k)
77         V = self.split_heads(V) # (batch, num_heads,
seq_len_k, d_k)
78
79         # Adjust mask for multiple heads
80         if mask is not None:
81             # Add head dimension: (batch, 1, seq_len_q,
seq_len_k)
82             mask = mask.unsqueeze(1)
83
84         # Scaled dot-product attention
85         d_k = Q.size(-1)
86         scores = torch.matmul(Q, K.transpose(-2, -1)) / math.
sqrt(d_k)
87
88         if mask is not None:
89             scores = scores.masked_fill(mask == 0, -1e9)
90
91         attention_weights = F.softmax(scores, dim=-1)
92         attention_weights = self.dropout(attention_weights)
93
94         # Apply attention to values
95         attention_output = torch.matmul(attention_weights, V)
96         # (batch, num_heads, seq_len_q, d_k)
97
98         # Combine heads
99         attention_output = self.combine_heads(
attention_output)
100         # (batch, seq_len_q, d_model)
101
102         # Final linear projection
103         output = self.W_o(attention_output)
104
105         return output, attention_weights
106
107     # Test
108     mha = MultiHeadAttention(d_model=512, num_heads=8)
109     x = torch.randn(2, 10, 512) # (batch, seq_len, d_model)
110     output, weights = mha(x, x, x) # Self-attention
111
112     print(f"Output shape: {output.shape}") # (2, 10, 512)

```

```
113 print(f"Attention weights shape: {weights.shape}") # (2, 8,  
10, 10)
```


11.4.3 Positional Encoding

```

1 class PositionalEncoding(nn.Module):
2     """Sinusoidal positional encoding."""
3
4     def __init__(self, d_model, max_len=5000, dropout=0.1):
5         super().__init__()
6         self.dropout = nn.Dropout(dropout)
7
8         # Create positional encoding matrix
9         pe = torch.zeros(max_len, d_model)
10        position = torch.arange(0, max_len, dtype=torch.float
11        ).unsqueeze(1)
12
13        # Compute div_term for scaling
14        div_term = torch.exp(
15            torch.arange(0, d_model, 2).float() *
16            (-math.log(10000.0) / d_model)
17        )
18
19        # Apply sine to even indices
20        pe[:, 0::2] = torch.sin(position * div_term)
21
22        # Apply cosine to odd indices
23        pe[:, 1::2] = torch.cos(position * div_term)
24
25        # Add batch dimension: (1, max_len, d_model)
26        pe = pe.unsqueeze(0)
27
28        # Register as buffer (not a parameter, but part of
29        state)
30        self.register_buffer('pe', pe)
31
32    def forward(self, x):
33        """
34        Args:
35        x: (batch, seq_len, d_model)
36        """
37        # Add positional encoding
38        x = x + self.pe[:, :x.size(1), :]
39        return self.dropout(x)
40
41    # Test
42    pe = PositionalEncoding(d_model=512, max_len=100)
43    x = torch.randn(2, 50, 512) # (batch, seq_len, d_model)
44    x_with_pos = pe(x)
45    print(x_with_pos.shape) # (2, 50, 512)
46
47    # Visualize positional encoding
48    import matplotlib.pyplot as plt
49    plt.figure(figsize=(12, 4))
50    plt.imshow(pe.pe[0, :100, :].detach().numpy(), aspect='auto',
51               cmap='viridis')
52    plt.xlabel('Dimension')
53    plt.ylabel('Position')
54    plt.colorbar()
55    plt.title('Positional Encoding')
56    plt.show()

```

11.4.4 Feed-Forward Network

```

1 class PositionWiseFeedForward(nn.Module):
2     """Position-wise feed-forward network."""
3
4     def __init__(self, d_model, d_ff, dropout=0.1):
5         """
6         Args:
7             d_model: Input/output dimension (e.g., 512)
8             d_ff: Hidden dimension (typically 4 * d_model =
2048)
9             dropout: Dropout probability
10        """
11        super().__init__()
12
13        self.fc1 = nn.Linear(d_model, d_ff)
14        self.fc2 = nn.Linear(d_ff, d_model)
15        self.dropout = nn.Dropout(dropout)
16
17    def forward(self, x):
18        """
19        Args:
20            x: (batch, seq_len, d_model)
21        """
22        # First linear + ReLU
23        x = F.relu(self.fc1(x))
24        x = self.dropout(x)
25
26        # Second linear
27        x = self.fc2(x)
28
29        return x
30
31    # Test
32    ffn = PositionWiseFeedForward(d_model=512, d_ff=2048)
33    x = torch.randn(2, 10, 512)
34    output = ffn(x)
35    print(output.shape)  # (2, 10, 512)

```

11.4.5 Complete Transformer Encoder Block

```

1 class TransformerEncoderBlock(nn.Module):
2     """Single Transformer encoder block."""
3
4     def __init__(self, d_model, num_heads, d_ff, dropout=0.1):
5         :
6         super().__init__()
7
8         # Multi-head attention
9         self.attention = MultiHeadAttention(d_model,
num_heads, dropout)
10
11        # Feed-forward network
12        self.ffn = PositionWiseFeedForward(d_model, d_ff,
dropout)
13
14        # Layer normalization
15        self.norm1 = nn.LayerNorm(d_model)
16        self.norm2 = nn.LayerNorm(d_model)
17
18        # Dropout

```

```

18     self.dropout1 = nn.Dropout(dropout)
19     self.dropout2 = nn.Dropout(dropout)
20
21     def forward(self, x, mask=None):
22         """
23         Args:
24             x: (batch, seq_len, d_model)
25             mask: Optional attention mask
26         """
27         # Multi-head self-attention with residual and layer
28         norm
29         attn_output, _ = self.attention(x, x, x, mask)
30         x = x + self.dropout1(attn_output) # Residual
31         connection
32         x = self.norm1(x) # Layer normalization
33
34         # Feed-forward with residual and layer norm
35         ffn_output = self.ffn(x)
36         x = x + self.dropout2(ffn_output) # Residual
37         connection
38         x = self.norm2(x) # Layer normalization
39
40         return x
41
42 # Test
43 encoder_block = TransformerEncoderBlock(d_model=512,
44                                         num_heads=8, d_ff=2048)
45 x = torch.randn(2, 10, 512)
46 output = encoder_block(x)
47 print(output.shape) # (2, 10, 512)

```

[Layer Norm Placement: Pre-Norm vs Post-Norm] **Original Transformer (post-norm):**

```
1 x = norm(x + sublayer(x))
```

Modern practice (pre-norm):

```
1 x = x + sublayer(norm(x))
```

Pre-norm is now preferred:

- More stable training
- Can train deeper networks
- Better gradient flow

The code above uses post-norm (original paper), but many modern implementations use pre-norm.

11.4.6 Using PyTorch's Built-in Attention

```

1  # PyTorch provides nn.MultiheadAttention
2  attention = nn.MultiheadAttention(
3      embed_dim=512,
4      num_heads=8,
5      dropout=0.1,
6      batch_first=True # Important! Use batch_first=True
7  )
8
9  # Self-attention
10 x = torch.randn(2, 10, 512) # (batch, seq_len, embed_dim)
11 attn_output, attn_weights = attention(x, x, x)
12
13 print(attn_output.shape) # (2, 10, 512)
14 print(attn_weights.shape) # (2, 10, 10) if need_weights=True
15
16 # Cross-attention (e.g., encoder-decoder)
17 encoder_output = torch.randn(2, 20, 512)
18 decoder_input = torch.randn(2, 10, 512)
19
20 attn_output, _ = attention(
21     query=decoder_input,
22     key=encoder_output,
23     value=encoder_output
24 )
25 print(attn_output.shape) # (2, 10, 512)

```

[batch_first in MultiheadAttention] PyTorch's nn.MultiheadAttention defaults to batch_first=False, expecting shape (seq_len, batch, embed_dim).

Always use batch_first=True for consistency:

```

1  # CONFUSING (default):
2  attention = nn.MultiheadAttention(embed_dim=512, num_heads=8)
3  x = torch.randn(10, 2, 512) # (seq_len, batch, embed_dim)
4
5  # CLEAR (recommended):
6  attention = nn.MultiheadAttention(embed_dim=512, num_heads=8,
7                                   batch_first=True)
8  x = torch.randn(2, 10, 512) # (batch, seq_len, embed_dim)

```

11.4.7 Masking for Causal Attention

For autoregressive models (e.g., language generation), prevent attending to future positions:

```

1  def create_causal_mask(seq_len):
2      """
3      Create causal mask to prevent attending to future
4      positions.
5
6      Returns:
7          mask: (seq_len, seq_len) with 1s on and below
8          diagonal
9      """
10     mask = torch.tril(torch.ones(seq_len, seq_len))
11     return mask
12
13 # Example
14 mask = create_causal_mask(5)

```

```

13 print(mask)
14 """
15 tensor([[1., 0., 0., 0., 0.],
16         [1., 1., 0., 0., 0.],
17         [1., 1., 1., 0., 0.],
18         [1., 1., 1., 1., 0.],
19         [1., 1., 1., 1., 1.]])
20 """
21
22 # Use with attention
23 seq_len = 10
24 causal_mask = create_causal_mask(seq_len)
25
26 # PyTorch MultiheadAttention expects mask where True = ignore
27 # So we invert: 0 = attend, 1 = ignore
28 attn_mask = (1 - causal_mask).bool()
29
30 x = torch.randn(2, seq_len, 512)
31 attention = nn.MultiheadAttention(embed_dim=512, num_heads=8,
32                                   batch_first=True)
33 output, _ = attention(x, x, x, attn_mask=attn_mask)

```

11.4.8 Complete Transformer for Classification

```

1 class TransformerClassifier(nn.Module):
2     """Transformer for sequence classification."""
3
4     def __init__(self, vocab_size, d_model, num_heads,
5                 num_layers,
6                 num_classes, d_ff, max_len=512, dropout=0.1)
7     :
8         super().__init__()
9
10        # Token embedding
11        self.embedding = nn.Embedding(vocab_size, d_model)
12
13        # Positional encoding
14        self.pos_encoding = PositionalEncoding(d_model,
15        max_len, dropout)
16
17        # Transformer encoder layers
18        encoder_layer = nn.TransformerEncoderLayer(
19            d_model=d_model,
20            nhead=num_heads,
21            dim_feedforward=d_ff,
22            dropout=dropout,
23            batch_first=True
24        )
25
26        self.transformer_encoder = nn.TransformerEncoder(
27            encoder_layer,
28            num_layers=num_layers
29        )
30
31        # Classification head
32        self.fc = nn.Linear(d_model, num_classes)
33
34        self.d_model = d_model
35
36        def forward(self, x, mask=None):
37            """

```

```

35     Args:
36         x: (batch, seq_len) - token indices
37         mask: Optional padding mask (batch, seq_len)
38     """
39     # Embedding and scaling
40     x = self.embedding(x) * math.sqrt(self.d_model)
41
42     # Add positional encoding
43     x = self.pos_encoding(x)
44
45     # Transformer encoding
46     x = self.transformer_encoder(x, src_key_padding_mask=
mask)
47
48     # Global average pooling
49     x = x.mean(dim=1) # (batch, d_model)
50
51     # Classification
52     logits = self.fc(x)
53
54     return logits
55
56 # Test
57 model = TransformerClassifier(
58     vocab_size=10000,
59     d_model=512,
60     num_heads=8,
61     num_layers=6,
62     num_classes=3,
63     d_ff=2048,
64     dropout=0.1
65 )
66
67 # Random token sequences
68 x = torch.randint(0, 10000, (32, 50)) # (batch, seq_len)
69 output = model(x)
70 print(output.shape) # (32, 3)
71
72 # Count parameters
73 total_params = sum(p.numel() for p in model.parameters())
74 print(f"Total parameters: {total_params:,}")

```

11.4.9 Transformer for Sequence-to-Sequence

```

1 class TransformerSeq2Seq(nn.Module):
2     """Complete Transformer with encoder and decoder."""
3
4     def __init__(self, src_vocab_size, tgt_vocab_size,
5                 d_model=512, num_heads=8, num_layers=6, d_ff=2048,
6                 dropout=0.1):
7         super().__init__()
8
9         # Source embeddings
10        self.src_embedding = nn.Embedding(src_vocab_size,
11                                         d_model)
12        self.src_pos_encoding = PositionalEncoding(d_model,
13                                                  dropout=dropout)
14
15        # Target embeddings
16        self.tgt_embedding = nn.Embedding(tgt_vocab_size,
17                                         d_model)
18        self.tgt_pos_encoding = PositionalEncoding(d_model,
19                                                  dropout=dropout)
20
21        # Transformer
22        self.transformer = nn.Transformer(
23            d_model=d_model,
24            nhead=num_heads,
25            num_encoder_layers=num_layers,
26            num_decoder_layers=num_layers,
27            dim_feedforward=d_ff,
28            dropout=dropout,
29            batch_first=True
30        )
31
32        # Output projection
33        self.fc_out = nn.Linear(d_model, tgt_vocab_size)
34
35        self.d_model = d_model
36
37    def forward(self, src, tgt, src_mask=None, tgt_mask=None,
38              src_padding_mask=None, tgt_padding_mask=None):
39        """
40        Args:
41        src: (batch, src_seq_len) - source token indices
42        tgt: (batch, tgt_seq_len) - target token indices
43        src_mask: Optional source attention mask
44        tgt_mask: Causal mask for target (prevent seeing
45        future)
46        src_padding_mask: Padding mask for source
47        tgt_padding_mask: Padding mask for target
48        """
49
50        # Embed and add positional encoding
51        src = self.src_embedding(src) * math.sqrt(self.
52        d_model)
53        src = self.src_pos_encoding(src)
54
55        tgt = self.tgt_embedding(tgt) * math.sqrt(self.
56        d_model)
57        tgt = self.tgt_pos_encoding(tgt)
58
59        # Transformer

```

```

51         output = self.transformer(
52             src, tgt,
53             src_mask=src_mask,
54             tgt_mask=tgt_mask,
55             src_key_padding_mask=src_padding_mask,
56             tgt_key_padding_mask=tgt_padding_mask
57         )
58
59         # Project to vocabulary
60         logits = self.fc_out(output)
61
62         return logits
63
64     def generate(self, src, max_len, start_token, end_token):
65         """
66         Generate sequence autoregressively.
67
68         Args:
69             src: (batch, src_seq_len)
70             max_len: Maximum generation length
71             start_token: Token to start generation
72             end_token: Token to stop generation
73         """
74         self.eval()
75         batch_size = src.size(0)
76         device = src.device
77
78         # Encode source
79         src = self.src_embedding(src) * math.sqrt(self.
d_model)
80         src = self.src_pos_encoding(src)
81         memory = self.transformer.encoder(src)
82
83         # Start with start_token
84         tgt = torch.full((batch_size, 1), start_token,
85                           dtype=torch.long, device=device)
86
87         for _ in range(max_len):
88             # Create causal mask
89             tgt_mask = nn.Transformer.
generate_square_subsequent_mask(
90                 tgt.size(1)
91             ).to(device)
92
93             # Decode
94             tgt_embedded = self.tgt_embedding(tgt) * math.
sqrt(self.d_model)
95             tgt_embedded = self.tgt_pos_encoding(tgt_embedded
)
96
97             output = self.transformer.decoder(tgt_embedded,
memory,
98                                               tgt_mask=
tgt_mask)
99
100             # Project to vocabulary
101             logits = self.fc_out(output[:, -1, :]) # Last
position
102
103             # Greedy decoding (take argmax)
104             next_token = logits.argmax(dim=-1, keepdim=True)
105

```



```

106         # Append to sequence
107         tgt = torch.cat([tgt, next_token], dim=1)
108
109         # Stop if all sequences generated end_token
110         if (next_token == end_token).all():
111             break
112
113     return tgt
114
115 # Test
116 model = TransformerSeq2Seq(
117     src_vocab_size=10000,
118     tgt_vocab_size=10000,
119     d_model=512,
120     num_heads=8,
121     num_layers=6
122 )
123
124 src = torch.randint(0, 10000, (2, 20)) # Source sequences
125 tgt = torch.randint(0, 10000, (2, 15)) # Target sequences
126
127 # Training forward pass
128 output = model(src, tgt[:, :-1]) # Shift target by 1
129 print(output.shape) # (2, 14, 10000)
130
131 # Generation
132 generated = model.generate(src, max_len=20, start_token=1,
133                             end_token=2)
134 print(generated.shape) # (2, <=20)

```

11.5 Exercises

11.1: Attention from Scratch - ★★

Goal: Implement and understand attention mechanism.

1. Implement scaled dot-product attention from scratch (without using the provided code)
2. Test on random Q, K, V matrices
3. Verify:
 - Attention weights sum to 1
 - Output shape is correct
 - Scaling by $\sqrt{d_k}$ stabilizes scores
4. Visualize attention weights as a heatmap

Visualization:

```

1 import matplotlib.pyplot as plt
2
3 # Compute attention
4 output, weights = scaled_dot_product_attention(Q, K, V)
5
6 # Visualize attention weights (first sample in batch)
7 plt.figure(figsize=(8, 6))
8 plt.imshow(weights[0].detach().numpy(), cmap='viridis')
9 plt.colorbar()
10 plt.xlabel('Key Position')
11 plt.ylabel('Query Position')
12 plt.title('Attention Weights')
13 plt.show()
```

11.2: Multi-Head Attention - ★★★

Goal: Implement multi-head attention.

1. Complete the MultiHeadAttention class
2. Test with different numbers of heads (1, 2, 4, 8)
3. Verify outputs are identical when using 1 head vs your scaled dot-product attention
4. Visualize attention patterns from different heads
5. Compare parameter counts: single-head vs multi-head

Questions:

- Do different heads learn different attention patterns?
- How does performance change with number of heads?

11.3: Positional Encoding Analysis - ★★★

Goal: Understand positional encoding.

1. Implement sinusoidal positional encoding
2. Visualize the encoding for positions 0-100
3. Experiment:
 - Train a Transformer with positional encoding
 - Train without positional encoding
 - Use learned positional embeddings instead

4. Compare:

- Final accuracy
- Does model learn order information without PE?
- Which PE works best?

Task: Sequence classification (e.g., sentiment analysis)

11.4: Complete Transformer Encoder - ★★★★★

Goal: Build a full Transformer encoder for classification.

1. Implement complete Transformer encoder:
 - Multi-head attention
 - Position-wise FFN
 - Layer normalization
 - Residual connections
 - Positional encoding
2. Train on text classification task (e.g., IMDB sentiment)
3. Use vocabulary size = 10,000 (most common words)
4. Hyperparameters:
 - $d_{model} = 256$
 - num_heads = 8
 - num_layers = 4
 - $d_{ff} = 1024$
5. Achieve >85% accuracy
6. Visualize attention weights on example sentences

Dataset preparation:

```

1 from torchtext.datasets import IMDB
2 from torchtext.data.utils import get_tokenizer
3 from collections import Counter
4
5 tokenizer = get_tokenizer('basic_english')
6
7 # Build vocabulary from training data
8 def build_vocab(data_iter, tokenizer, max_size=10000):
9     counter = Counter()
10    for label, text in data_iter:
11        counter.update(tokenizer(text))
12
13    vocab = {word: i+2 for i, (word, _) in
14            enumerate(counter.most_common(max_size))}
15    vocab['<pad>'] = 0
16    vocab['<unk>'] = 1
17
18    return vocab
19
20 # Encode text
21 def encode_text(text, vocab, tokenizer, max_len=256):
22     tokens = tokenizer(text)[:max_len]
23     indices = [vocab.get(token, vocab['<unk>']) for
24               token in tokens]
25
26     # Pad to max_len
27     indices += [vocab['<pad>']] * (max_len - len(indices))

```

```

27     ))
28     return torch.LongTensor(indices)

```

11.5: Causal Attention for Language Modeling - ★★★★★

Goal: Implement autoregressive generation with causal masking.

1. Build a decoder-only Transformer (like GPT)
2. Use causal masking (can't attend to future tokens)
3. Train on character-level language modeling:
 - Dataset: Shakespeare text or similar
 - Task: Predict next character
4. Implement generation:
 - Greedy decoding (argmax)
 - Top-k sampling
 - Temperature sampling
5. Generate text samples and evaluate quality

Generation strategies:

```

1 def generate_text(model, start_text, max_len=100,
2                   temperature=1.0,
3                   top_k=None):
4     """Generate text autoregressively."""
5     model.eval()
6
7     # Encode start text
8     tokens = encode(start_text)
9
10    with torch.no_grad():
11        for _ in range(max_len):
12            # Forward pass
13            logits = model(tokens)
14
15            # Get logits for last position
16            logits = logits[-1] / temperature
17
18            # Top-k filtering
19            if top_k is not None:
20                values, _ = torch.topk(logits, top_k)
21                logits[logits < values[-1]] = -float('
22                inf')
23
24            # Sample from distribution
25            probs = F.softmax(logits, dim=-1)
26            next_token = torch.multinomial(probs, 1)
27
28            # Append to sequence
29            tokens = torch.cat([tokens, next_token])
30
31    return decode(tokens)

```

11.6: Attention Visualization - ★★★★★

Goal: Understand what Transformers learn.

1. Train a Transformer on sequence task
2. Extract attention weights from different layers and heads
3. Visualize:
 - Which heads focus on local context?
 - Which heads capture long-range dependencies?
 - How do patterns differ across layers?
4. Create attention heatmaps for example sentences
5. Analyze:
 - Do lower layers learn syntax?
 - Do higher layers learn semantics?

Visualization code:

```

1 def visualize_attention_heads(model, sentence, tokenizer
2 ):
3     """Visualize attention patterns from all heads."""
4
5     # Encode sentence
6     tokens = encode(sentence, tokenizer)
7
8     # Forward pass with hooks to capture attention
9     attentions = []
10
11     def hook_fn(module, input, output):
12         # output[1] is attention weights
13         attentions.append(output[1].detach())
14
15     hooks = []
16     for layer in model.transformer_encoder.layers:
17         hook = layer.self_attn.register_forward_hook(
18             hook_fn)
19         hooks.append(hook)
20
21     # Forward
22     _ = model(tokens.unsqueeze(0))
23
24     # Remove hooks
25     for hook in hooks:
26         hook.remove()
27
28     # Plot attention from each layer and head
29     num_layers = len(attentions)
30     num_heads = attentions[0].size(1)
31
32     fig, axes = plt.subplots(num_layers, num_heads,
33                             figsize=(num_heads*3,
34                                     num_layers*3))
35
36     for layer in range(num_layers):
37         for head in range(num_heads):
38             ax = axes[layer, head]
39             attn = attentions[layer][0, head].cpu().
40             numpy()
41
42             im = ax.imshow(attn, cmap='viridis')
43             ax.set_title(f'L{layer+1} H{head+1}')
```

```
40         ax.set_xlabel('Key')
41         ax.set_ylabel('Query')
42
43     plt.tight_layout()
44     plt.show()
45
46     # Example
47     visualize_attention_heads(model, "The cat sat on the mat
    ", tokenizer)
```

11.6 Key Takeaways

Why Transformers dominate:

- **Parallelization:** All positions processed simultaneously (10-100× faster than RNNs)
- **Long-range dependencies:** Direct connections between any positions
- **Scalability:** Works with billions of parameters and tokens
- **Flexibility:** Minimal inductive biases, learns from data

Core components:

- **Scaled dot-product attention:** $\text{Attention}(Q, K, V) = \text{softmax}(QK^T/\sqrt{d_k})V$
- **Multi-head attention:** Multiple attention mechanisms in parallel
- **Positional encoding:** Inject position information (no inherent order)
- **Feed-forward networks:** Position-wise transformation
- **Residual connections + Layer norm:** Stable training

Architecture patterns:

- **Encoder-only:** BERT (classification, understanding)
- **Decoder-only:** GPT (generation, autoregressive)
- **Encoder-decoder:** T5, BART (translation, summarization)

Implementation tips:

- Always use `batch_first=True` in PyTorch
- Scale attention scores by $\sqrt{d_k}$ (stability)
- Use causal masking for autoregressive tasks
- Pre-norm is more stable than post-norm for deep networks
- Typical hyperparameters: $d_{model} = 512$, $h = 8$, $d_{ff} = 2048$

Training considerations:

- Transformers need lots of data (millions of examples)
- Use warmup learning rate schedule
- Gradient clipping less critical than RNNs (but still useful)
- Dropout on attention weights (0.1) and residual connections
- Label smoothing helps (0.1)

When to use Transformers:

- NLP tasks (now standard)
- Long sequences (>100 tokens)
- When you have lots of data
- When parallelization matters
- Vision tasks (Vision Transformer - ViT)

Limitations:

- **Quadratic complexity:** $O(n^2)$ in sequence length (memory and compute)
- **Data hungry:** Requires large datasets for good performance
- **Memory intensive:** Attention matrices can be huge
- **Position limit:** Most models have fixed maximum sequence length

Solutions to limitations:

- Sparse attention (only attend to subset of positions)
- Linear attention approximations (e.g., Linformer)
- Local attention windows (e.g., Longformer)
- Efficient Transformers (e.g., Reformer, Performer)

Common mistakes:

- Forgetting positional encoding (model can't distinguish positions)
- Wrong mask shape or type
- Not using causal mask for generation
- Forgetting to scale attention scores
- Using post-norm for very deep networks (use pre-norm)

Modern developments:

- **Vision Transformers (ViT):** Transformers for images
- **BERT:** Bidirectional encoder (masked language modeling)
- **GPT:** Decoder-only (autoregressive generation)
- **T5:** Unified text-to-text framework
- **LLaMA, ChatGPT:** Large language models at scale

Part III

Practical Skills

12 Debugging & Best Practices

12.1 Introduction: Why Debugging Matters

Deep learning is full of silent failures:

- Model trains but doesn't learn (loss stuck at baseline)
- Loss is NaN after a few iterations
- Training loss decreases but validation doesn't
- Model works on toy data but fails on real data
- Gradients vanish or explode

Unlike traditional programming where bugs crash immediately, deep learning bugs often manifest as **poor performance**. You need systematic debugging strategies.

12.2 Common Bugs and How to Fix Them

12.2.1 Bug 1: Loss is NaN

Symptoms:

```
1 Epoch 1: loss = 2.345
2 Epoch 2: loss = 1.876
3 Epoch 3: loss = nan
```

Causes and fixes:

1. Learning rate too high

```
1 # Check: Does loss explode before NaN?
2 # Fix: Reduce learning rate
3 optimizer = torch.optim.Adam(model.parameters(), lr=1e-5) #
    Try lower
```

2. Numerical instability in loss

```
1 # BAD: Manual log can cause NaN
2 loss = -torch.log(predictions) # NaN if predictions = 0
3
4 # GOOD: Use stable implementations
5 loss = F.cross_entropy(logits, targets) # Numerically stable
6 # or
7 loss = F.binary_cross_entropy_with_logits(logits, targets)
```

3. Exploding gradients (RNNs)

```
1 # Fix: Add gradient clipping
2 torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm
    =1.0)
```

4. Data contains NaN or Inf

```

1 # Check data
2 print(torch.isnan(data).any())
3 print(torch.isinf(data).any())
4
5 # Fix: Clean data or add checks
6 data = torch.nan_to_num(data, nan=0.0, posinf=1e6, neginf=-1
    e6)

```

Detection code:

```

1 def check_for_nan(model, data, target):
2     """Debug NaN issues."""
3     # Check inputs
4     print(f"Data has NaN: {torch.isnan(data).any()}")
5     print(f"Data has Inf: {torch.isinf(data).any()}")
6
7     # Forward pass
8     output = model(data)
9     print(f"Output has NaN: {torch.isnan(output).any()}")
10
11    # Backward pass
12    loss = criterion(output, target)
13    print(f"Loss has NaN: {torch.isnan(loss).any()}")
14    loss.backward()
15
16    # Check gradients
17    for name, param in model.named_parameters():
18        if param.grad is not None:
19            if torch.isnan(param.grad).any():
20                print(f"NaN gradient in {name}")
21            if torch.isinf(param.grad).any():
22                print(f"Inf gradient in {name}")

```

12.2.2 Bug 2: Loss Not Decreasing

Symptoms:

```

1 Epoch 1: loss = 2.345
2 Epoch 2: loss = 2.342
3 Epoch 3: loss = 2.340
4 ...
5 Epoch 100: loss = 2.300  # Barely moving

```

Causes and fixes:

1. Learning rate too low

```

1 # Fix: Increase learning rate
2 optimizer = torch.optim.Adam(model.parameters(), lr=1e-2) #
   Try higher
3
4 # Or use learning rate finder
5 def find_lr(model, train_loader, init_lr=1e-8, final_lr=10):
6     """Find optimal learning rate."""
7     lrs = []
8     losses = []
9
10    optimizer = torch.optim.SGD(model.parameters(), lr=
init_lr)
11    lr_scheduler = torch.optim.lr_scheduler.ExponentialLR(
12        optimizer, gamma=(final_lr/init_lr)**(1/100)
13    )
14
15    for batch in train_loader:
16        optimizer.zero_grad()
17        loss = compute_loss(model, batch)
18        loss.backward()
19        optimizer.step()
20
21        lrs.append(optimizer.param_groups[0]['lr'])
22        losses.append(loss.item())
23
24        lr_scheduler.step()
25
26        if len(lrs) >= 100:
27            break
28
29    # Plot
30    import matplotlib.pyplot as plt
31    plt.plot(lrs, losses)
32    plt.xscale('log')
33    plt.xlabel('Learning Rate')
34    plt.ylabel('Loss')
35    plt.title('Learning Rate Finder')
36    plt.show()
37
38    # Optimal LR is usually where loss decreases fastest
39    # (steepest negative gradient)

```

2. Wrong loss function

```

1 # BAD: Using MSE for classification
2 criterion = nn.MSELoss() # Wrong for discrete labels!
3
4 # GOOD: Use appropriate loss

```

```
5 criterion = nn.CrossEntropyLoss() # For classification
```

3. Forgot to zero gradients

```
1 # BAD: Gradients accumulate
2 for batch in dataloader:
3     loss = criterion(model(x), y)
4     loss.backward()
5     optimizer.step() # FORGOT optimizer.zero_grad()!
6
7 # GOOD:
8 for batch in dataloader:
9     optimizer.zero_grad() # Always zero first!
10    loss = criterion(model(x), y)
11    loss.backward()
12    optimizer.step()
```

4. Data not normalized

```
1 # Check data range
2 print(f"Data min: {data.min()}, max: {data.max()}")
3
4 # Fix: Normalize
5 mean = data.mean()
6 std = data.std()
7 data = (data - mean) / (std + 1e-8)
```

5. Dead ReLUs or vanishing gradients

```
1 # Check activation distributions
2 def check_activations(model, data):
3     activations = {}
4
5     def hook(name):
6         def hook_fn(module, input, output):
7             activations[name] = output.detach()
8             return hook_fn
9
10    # Register hooks
11    for name, module in model.named_modules():
12        if isinstance(module, nn.ReLU):
13            module.register_forward_hook(hook(name))
14
15    model(data)
16
17    # Check what percentage is zero
18    for name, act in activations.items():
19        zero_pct = (act == 0).float().mean() * 100
20        print(f"{name}: {zero_pct:.1f}% zeros")
21
22        if zero_pct > 50:
23            print(f" WARNING: More than 50% dead neurons!")
24
25    # Fix: Use Leaky ReLU or adjust initialization
```

12.2.3 Bug 3: Training Loss Decreases, Validation Doesn't

Symptoms: Overfitting

Causes and fixes:

1. Model too complex for data

```

1 # Fix: Reduce model size
2 model = SmallNet() # Fewer layers, fewer parameters
3
4 # Or increase data
5 # Or add regularization (see below)

```

2. No regularization

```

1 # Add dropout
2 model = nn.Sequential(
3     nn.Linear(100, 128),
4     nn.ReLU(),
5     nn.Dropout(0.5), # Add dropout
6     nn.Linear(128, 10)
7 )
8
9 # Add weight decay
10 optimizer = torch.optim.Adam(model.parameters(), lr=1e-3,
11                               weight_decay=1e-4) # L2
12                               regularization
13
14 # Add data augmentation
15 transform = transforms.Compose([
16     transforms.RandomHorizontalFlip(),
17     transforms.RandomCrop(32, padding=4),
18     transforms.ToTensor()
19 ])

```

3. Training too long

```

1 # Implement early stopping
2 class EarlyStopping:
3     def __init__(self, patience=10, min_delta=0):
4         self.patience = patience
5         self.min_delta = min_delta
6         self.counter = 0
7         self.best_loss = None
8         self.should_stop = False
9
10    def __call__(self, val_loss):
11        if self.best_loss is None:
12            self.best_loss = val_loss
13        elif val_loss > self.best_loss - self.min_delta:
14            self.counter += 1
15            if self.counter >= self.patience:
16                self.should_stop = True
17        else:
18            self.best_loss = val_loss
19            self.counter = 0
20
21 # Usage
22 early_stopping = EarlyStopping(patience=10)
23
24 for epoch in range(max_epochs):

```

```

25     train_loss = train_epoch()
26     val_loss = validate()
27
28     early_stopping(val_loss)
29     if early_stopping.should_stop:
30         print(f"Early stopping at epoch {epoch}")
31         break

```

12.2.4 Bug 4: Model Works on Small Data, Fails on Full Dataset

Causes and fixes:

1. Batch size too small

```

1  # Batch norm unstable with small batches
2  # Fix: Increase batch size or use Layer Norm
3
4  # Or use Gradient Accumulation
5  accumulation_steps = 4
6
7  for i, batch in enumerate(dataloader):
8      loss = criterion(model(x), y)
9      loss = loss / accumulation_steps # Scale loss
10     loss.backward()
11
12     if (i + 1) % accumulation_steps == 0:
13         optimizer.step()
14         optimizer.zero_grad()

```

2. Data distribution shift

```

1  # Check: Are train and val distributions similar?
2  def check_distribution(train_data, val_data):
3      print(f"Train mean: {train_data.mean()}, std: {train_data.std()}")
4      print(f"Val mean: {val_data.mean()}, std: {val_data.std()}")
5
6      # Visualize
7      import matplotlib.pyplot as plt
8      plt.hist(train_data.flatten(), bins=50, alpha=0.5, label='Train')
9      plt.hist(val_data.flatten(), bins=50, alpha=0.5, label='Val')
10     plt.legend()
11     plt.show()
12
13 # Fix: Normalize using training statistics
14 mean = train_data.mean()
15 std = train_data.std()
16
17 train_data = (train_data - mean) / std
18 val_data = (val_data - mean) / std # Use training stats!

```

12.3 Debugging Strategies

12.3.1 Start Simple, Add Complexity

Step 1: Overfit on one batch

If you can't overfit one batch, something is fundamentally wrong.

```

1  # Get one batch
2  x, y = next(iter(train_loader))
3
4  # Try to overfit it
5  for epoch in range(1000):
6      optimizer.zero_grad()
7      output = model(x)
8      loss = criterion(output, y)
9      loss.backward()
10     optimizer.step()
11
12     if epoch % 100 == 0:
13         print(f"Epoch {epoch}: Loss = {loss.item():.4f}")
14
15 # Should reach near-zero loss
16 # If not, check:
17 # - Model has enough capacity
18 # - Loss function is correct
19 # - Learning rate is reasonable

```

Step 2: Overfit on small dataset

Try 100 samples. Should still overfit easily.

Step 3: Add validation

Once overfitting works, add validation and regularization.

Step 4: Scale to full dataset

Now train on full data with proper regularization.

12.3.2 Check Shapes at Every Step

```

1  def debug_shapes(model, x):
2      """Print shapes through the network."""
3      print(f"Input: {x.shape}")
4
5      for i, (name, module) in enumerate(model.named_children()):
6          x = module(x)
7          print(f"{i}. {name}: {x.shape}")
8
9      return x
10
11 # Usage
12 x = torch.randn(4, 3, 32, 32)
13 output = debug_shapes(model, x)
14 """
15 Input: torch.Size([4, 3, 32, 32])
16 0. conv1: torch.Size([4, 64, 32, 32])
17 1. relu: torch.Size([4, 64, 32, 32])
18 2. pool: torch.Size([4, 64, 16, 16])
19 ...
20 """

```

12.3.3 Visualize Gradients

```

1 def plot_grad_flow(named_parameters):
2     """
3     Plots gradient flow through the network.
4     Helps identify vanishing/exploding gradients.
5     """
6     ave_grads = []
7     max_grads = []
8     layers = []
9
10    for n, p in named_parameters:
11        if p.requires_grad and p.grad is not None:
12            layers.append(n)
13            ave_grads.append(p.grad.abs().mean().item())
14            max_grads.append(p.grad.abs().max().item())
15
16    import matplotlib.pyplot as plt
17    plt.figure(figsize=(12, 6))
18    plt.bar(range(len(ave_grads)), ave_grads, alpha=0.5,
19            label='mean')
20    plt.bar(range(len(max_grads)), max_grads, alpha=0.5,
21            label='max')
22    plt.hlines(0, 0, len(ave_grads), linewidth=2, color='k')
23    plt.xticks(range(len(layers)), layers, rotation='vertical')
24    plt.xlabel('Layers')
25    plt.ylabel('Gradient magnitude')
26    plt.legend()
27    plt.title('Gradient Flow')
28    plt.grid(True)
29    plt.tight_layout()
30    plt.show()
31
32    # Usage after loss.backward()
33    plot_grad_flow(model.named_parameters())

```

12.3.4 Use Hooks for Debugging

```

1 # Forward hook: inspect activations
2 def forward_hook(module, input, output):
3     print(f"Forward: {output.shape}, mean={output.mean():.4f}")
4
5 # Backward hook: inspect gradients
6 def backward_hook(module, grad_input, grad_output):
7     print(f"Backward: grad mean={grad_output[0].mean():.4f}")
8
9 # Register hooks
10 for name, module in model.named_modules():
11     if isinstance(module, nn.Linear):
12         module.register_forward_hook(forward_hook)
13         module.register_backward_hook(backward_hook)
14
15 # Now run forward and backward
16 output = model(x)
17 loss = criterion(output, y)
18 loss.backward()
19
20 # Hooks will print debug info

```


12.4 Best Practices for Development

12.4.1 Code Organization

Project structure:

```
project/
  data/
    __init__.py
    dataset.py          # Dataset classes
    transforms.py       # Data augmentation
  models/
    __init__.py
    resnet.py           # Model architectures
    utils.py            # Model utilities
  utils/
    __init__.py
    train.py            # Training utilities
    eval.py             # Evaluation utilities
    metrics.py          # Custom metrics
  configs/
    config.yaml         # Hyperparameters
  train.py              # Main training script
  evaluate.py           # Evaluation script
  requirements.txt      # Dependencies
```

Separate concerns:

```
1  # models/resnet.py
2  class ResNet(nn.Module):
3      """Model definition only."""
4      def __init__(self, ...):
5          # Architecture
6
7      def forward(self, x):
8          # Forward pass
9          return x
10
11 # utils/train.py
12 def train_epoch(model, loader, optimizer, criterion, device):
13     """Training logic."""
14     model.train()
15     total_loss = 0
16
17     for batch in loader:
18         # Training step
19         ...
20
21     return total_loss / len(loader)
22
23 # train.py
24 def main():
25     # Configuration
26     # Data loading
27     # Model creation
28     # Training loop
29     pass
30
31 if __name__ == "__main__":
32     main()
```

12.4.2 Configuration Management

```

1  # config.yaml
2  model:
3      type: resnet18
4      num_classes: 10
5
6  training:
7      batch_size: 128
8      epochs: 100
9      learning_rate: 0.001
10     weight_decay: 0.0001
11
12  data:
13     root: ./data
14     num_workers: 4
15     augmentation: true
16
17  # Load config
18  import yaml
19
20  with open('config.yaml', 'r') as f:
21     config = yaml.safe_load(f)
22
23  # Use config
24  batch_size = config['training']['batch_size']
25  lr = config['training']['learning_rate']

```

12.4.3 Experiment Tracking

```

1  # Simple logging
2  import logging
3
4  logging.basicConfig(
5      level=logging.INFO,
6      format='%(asctime)s - %(levelname)s - %(message)s',
7      handlers=[
8          logging.FileHandler('training.log'),
9          logging.StreamHandler()
10     ]
11 )
12
13 logger = logging.getLogger(__name__)
14
15 # Use in training
16 logger.info(f"Epoch {epoch}: Train Loss = {train_loss:.4f}")
17
18 # TensorBoard (better)
19 from torch.utils.tensorboard import SummaryWriter
20
21 writer = SummaryWriter('runs/experiment_1')
22
23 for epoch in range(num_epochs):
24     train_loss = train_epoch()
25     val_loss = validate()
26
27     writer.add_scalar('Loss/train', train_loss, epoch)
28     writer.add_scalar('Loss/val', val_loss, epoch)
29     writer.add_scalar('LR', optimizer.param_groups[0]['lr'],
30                        epoch)

```

```
31 writer.close()  
32  
33 # View: tensorboard --logdir=runs
```

12.4.4 Reproducibility

```

1 import random
2 import numpy as np
3 import torch
4
5 def set_seed(seed=42):
6     """Set all random seeds for reproducibility."""
7     random.seed(seed)
8     np.random.seed(seed)
9     torch.manual_seed(seed)
10    torch.cuda.manual_seed(seed)
11    torch.cuda.manual_seed_all(seed)
12
13    # Make cudnn deterministic
14    torch.backends.cudnn.deterministic = True
15    torch.backends.cudnn.benchmark = False
16
17    # Call at start of training
18    set_seed(42)
19
20    # Note: This makes training slower but reproducible
21    # For production: set benchmark=True for speed

```

12.4.5 Checkpointing

```

1 def save_checkpoint(model, optimizer, epoch, loss, path):
2     """Save model checkpoint."""
3     torch.save({
4         'epoch': epoch,
5         'model_state_dict': model.state_dict(),
6         'optimizer_state_dict': optimizer.state_dict(),
7         'loss': loss,
8     }, path)
9
10    def load_checkpoint(model, optimizer, path):
11        """Load model checkpoint."""
12        checkpoint = torch.load(path)
13        model.load_state_dict(checkpoint['model_state_dict'])
14        optimizer.load_state_dict(checkpoint['
15        optimizer_state_dict'])
16        epoch = checkpoint['epoch']
17        loss = checkpoint['loss']
18        return epoch, loss
19
20    # Save best model
21    best_loss = float('inf')
22
23    for epoch in range(num_epochs):
24        train_loss = train_epoch()
25        val_loss = validate()
26
27        # Save if best
28        if val_loss < best_loss:
29            best_loss = val_loss
30            save_checkpoint(model, optimizer, epoch, val_loss,
31                            'best_model.pth')
32
33        # Save periodic checkpoints
34        if epoch % 10 == 0:
35            save_checkpoint(model, optimizer, epoch, val_loss,

```

12.5 Debugging Checklist

When starting a new project:

1. Can model overfit one batch?
2. Can model overfit small dataset (100 samples)?
3. Are shapes correct at every layer?
4. Is loss function appropriate for task?
5. Is learning rate reasonable? (Try LR finder)
6. Is data normalized?
7. Are you using the right optimizer?
8. Did you set `model.train()` / `model.eval()`?

When loss is NaN:

1. Check for NaN in data
2. Reduce learning rate
3. Add gradient clipping
4. Use stable loss functions
5. Check for division by zero

When loss doesn't decrease:

1. Try higher learning rate
2. Check if you're zeroing gradients
3. Verify loss function is correct
4. Check data normalization
5. Verify model has enough capacity
6. Check for dead ReLUs

When overfitting:

1. Add dropout
2. Add weight decay
3. Use data augmentation
4. Reduce model size
5. Get more data
6. Use early stopping

12.6 Key Takeaways

Systematic debugging:

- Start simple (overfit one batch)
- Add complexity gradually
- Check shapes at every step
- Visualize gradients and activations
- Use hooks for inspection

Common bugs:

- NaN loss: LR too high, numerical instability, exploding gradients
- Loss not decreasing: LR too low, wrong loss, forgot `zero_grad`, bad data
- Overfitting: Model too complex, no regularization, training too long
- Works on toy data, fails on real: Batch size, distribution shift

Best practices:

- Separate code into modules (data, models, training)
- Use configuration files
- Track experiments (TensorBoard)
- Set random seeds for reproducibility
- Save checkpoints regularly
- Log everything

Development workflow:

1. Start with simple baseline
2. Verify can overfit small data
3. Add regularization
4. Scale to full dataset
5. Tune hyperparameters
6. Validate on held-out test set

Tools for debugging:

- Print statements (shapes, values, statistics)
- Hooks (forward and backward)
- TensorBoard (visualize training)
- Gradient visualization
- Learning rate finder
- Activation distribution checks

Remember:

- Deep learning bugs are often silent (poor performance, not crashes)
- Always start with simplest possible setup
- Check your data first (garbage in = garbage out)
- One change at a time (isolate what works)
- Keep good records (what did you try, what worked)

13 Training Best Practices

13.1 Introduction: From Working to Working Well

You've built a model that trains. Now: how do you make it train **well**?

This section covers practical strategies for:

- Choosing hyperparameters
- Learning rate strategies
- Efficient training
- Avoiding common pitfalls

13.2 Hyperparameter Tuning

13.2.1 Which Hyperparameters Matter Most?

Priority order:

1. Learning rate (MOST IMPORTANT)

Can make 10× difference in performance. Get this right first.

2. Batch size

Affects both training speed and generalization.

3. Architecture (width, depth)

Number of layers and neurons per layer.

4. Regularization (dropout, weight decay)

Prevents overfitting.

5. Learning rate schedule

How LR changes over time.

6. Everything else

Optimizer choice, activation functions, initialization—usually matter less.

13.2.2 Learning Rate Strategies

Strategy 1: Learning Rate Finder

Find optimal LR before training:

```
1 from torch_lr_finder import LRFinder
2
3 model = YourModel()
4 optimizer = torch.optim.Adam(model.parameters(), lr=1e-7)
5 criterion = nn.CrossEntropyLoss()
6
7 lr_finder = LRFinder(model, optimizer, criterion, device="
  cuda")
8 lr_finder.range_test(train_loader, end_lr=10, num_iter=100)
9 lr_finder.plot() # Shows loss vs LR
10
11 # Pick LR where loss decreases fastest (steepest slope)
12 # Usually 10 smaller than where loss explodes
13 optimal_lr = 1e-3 # Read from plot
14
```

```
15 lr_finder.reset() # Reset model and optimizer
```

Strategy 2: Start Large, Decay

```
1 # Common schedules
2
3 # 1. Step decay (reduce every N epochs)
4 scheduler = torch.optim.lr_scheduler.StepLR(
5     optimizer, step_size=30, gamma=0.1
6 )
7 # LR    0.1 every 30 epochs
8
9 # 2. Exponential decay
10 scheduler = torch.optim.lr_scheduler.ExponentialLR(
11     optimizer, gamma=0.95
12 )
13 # LR    0.95 every epoch
14
15 # 3. Cosine annealing (smooth decay)
16 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(
17     optimizer, T_max=100, eta_min=1e-6
18 )
19 # Smooth cosine decay from initial_lr to eta_min over T_max
    epochs
20
21 # 4. Reduce on plateau (adaptive)
22 scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(
23     optimizer, mode='min', factor=0.5, patience=10
24 )
25 # Reduce LR by 0.5 if val loss doesn't improve for 10 epochs
26
27 # Usage
28 for epoch in range(num_epochs):
29     train_loss = train_epoch()
30     val_loss = validate()
31
32     # For ReduceLROnPlateau, pass validation loss
33     scheduler.step(val_loss) # or scheduler.step() for
        others
```

Strategy 3: Warmup + Cosine Decay

State-of-the-art for Transformers:

```
1 from torch.optim.lr_scheduler import LambdaLR
2
3 def get_cosine_schedule_with_warmup(optimizer,
4     num_warmup_steps,
5     num_training_steps):
6     """Warmup + cosine decay."""
7
8     def lr_lambda(current_step):
9         # Warmup
10         if current_step < num_warmup_steps:
11             return float(current_step) / float(max(1,
12                 num_warmup_steps))
13
14         # Cosine decay
15         progress = float(current_step - num_warmup_steps) / \
16             float(max(1, num_training_steps -
17                 num_warmup_steps))
```



```
15         return max(0.0, 0.5 * (1.0 + math.cos(math.pi *
16             progress)))
17     return LambdaLR(optimizer, lr_lambda)
18
19 # Usage
20 num_epochs = 100
21 steps_per_epoch = len(train_loader)
22 num_training_steps = num_epochs * steps_per_epoch
23 num_warmup_steps = 5 * steps_per_epoch # 5 epochs warmup
24
25 optimizer = torch.optim.AdamW(model.parameters(), lr=1e-3)
26 scheduler = get_cosine_schedule_with_warmup(
27     optimizer, num_warmup_steps, num_training_steps
28 )
29
30 for epoch in range(num_epochs):
31     for batch in train_loader:
32         # Training step
33         ...
34         optimizer.step()
35         scheduler.step() # Update LR every step, not epoch!
```

13.2.3 Batch Size Selection

Trade-offs:

Batch Size	Pros	Cons
Small (8-32)	Better generalization Less memory	Slower, noisy gradients Batch norm unstable
Medium (64-256)	Good balance	Standard choice
Large (512+)	Faster training More stable	Worse generalization Requires careful tuning

Table 9: Batch size trade-offs

Rules of thumb:

- **Start with 32-128:** Usually works well
- **Larger batches need higher LR:** Scale by $\sqrt{\text{batch size}}$
- Use gradient accumulation if memory limited
- Don't go below 8 if using batch norm

Linear scaling rule:

If you double batch size, double learning rate (approximately).

```

1 # Base setup
2 batch_size = 64
3 lr = 0.001
4
5 # Double batch size
6 batch_size = 128
7 lr = 0.002 # Double LR
8
9 # But: this is approximate! Always validate on your task

```

13.2.4 Model Architecture Selection

Width vs Depth:

```

1 # Wider networks (fewer layers, more neurons)
2 model_wide = MLP([100, 1024, 1024, 10]) # 2 hidden layers,
    1024 neurons
3
4 # Deeper networks (more layers, fewer neurons)
5 model_deep = MLP([100, 256, 256, 256, 256, 10]) # 4 hidden
    layers, 256 neurons
6
7 # Modern preference: deeper is usually better (with skip
    connections)

```

Start small, scale up:

1. Start with small model that trains fast
2. Check if it underfits (high train loss)
3. If underfitting: increase capacity
4. If overfitting: add regularization

13.3 Regularization Strategies

13.3.1 Dropout

```

1 # Standard dropout
2 model = nn.Sequential(
3     nn.Linear(100, 256),
4     nn.ReLU(),
5     nn.Dropout(0.5), # Drop 50% during training
6     nn.Linear(256, 10)
7 )
8
9 # Guidelines:
10 # - Start with 0.5, tune if needed
11 # - Higher dropout = more regularization
12 # - Use lower dropout (0.1-0.2) for convolutional layers
13 # - Higher dropout (0.5-0.7) for large fully connected layers

```

13.3.2 Weight Decay (L2 Regularization)

```

1 # Add weight decay to optimizer
2 optimizer = torch.optim.Adam(
3     model.parameters(),
4     lr=1e-3,
5     weight_decay=1e-4 # L2 penalty
6 )
7
8 # Typical values: 1e-5 to 1e-3
9 # Larger weight decay = stronger regularization
10
11 # Note: Use AdamW for better weight decay
12 optimizer = torch.optim.AdamW(model.parameters(), lr=1e-3,
13                                weight_decay=1e-4)

```

13.3.3 Data Augmentation

```

1 # For images
2 train_transform = transforms.Compose([
3     transforms.RandomHorizontalFlip(),
4     transforms.RandomCrop(32, padding=4),
5     transforms.ColorJitter(brightness=0.2, contrast=0.2),
6     transforms.RandomRotation(15),
7     transforms.ToTensor(),
8     transforms.Normalize(mean=[0.485, 0.456, 0.406],
9                             std=[0.229, 0.224, 0.225])
10 ])
11
12 # NO augmentation for validation/test
13 val_transform = transforms.Compose([
14     transforms.ToTensor(),
15     transforms.Normalize(mean=[0.485, 0.456, 0.406],
16                             std=[0.229, 0.224, 0.225])
17 ])

```

13.4 Training Monitoring

13.4.1 What to Track

Essential metrics:

- Training loss (per epoch)
- Validation loss (per epoch)
- Validation accuracy/metric (per epoch)
- Learning rate (per epoch or step)

Useful but optional:

- Gradient norms
- Weight norms
- Training time per epoch
- GPU memory usage

```

1  # Complete training loop with monitoring
2  from torch.utils.tensorboard import SummaryWriter
3
4  writer = SummaryWriter('runs/experiment')
5
6  for epoch in range(num_epochs):
7      # Training
8      model.train()
9      train_loss = 0
10     train_correct = 0
11     train_total = 0
12
13     for batch in train_loader:
14         optimizer.zero_grad()
15
16         x, y = batch
17         x, y = x.to(device), y.to(device)
18
19         output = model(x)
20         loss = criterion(output, y)
21         loss.backward()
22
23         # Track gradient norm
24         grad_norm = torch.nn.utils.clip_grad_norm_(
25             model.parameters(), max_norm=1.0
26         )
27
28         optimizer.step()
29
30         train_loss += loss.item()
31         pred = output.argmax(dim=1)
32         train_correct += (pred == y).sum().item()
33         train_total += y.size(0)
34
35     # Validation
36     model.eval()
37     val_loss = 0
38     val_correct = 0
39     val_total = 0
40
41     with torch.no_grad():
42         for batch in val_loader:
43             x, y = batch

```

```

44         x, y = x.to(device), y.to(device)
45
46         output = model(x)
47         loss = criterion(output, y)
48
49         val_loss += loss.item()
50         pred = output.argmax(dim=1)
51         val_correct += (pred == y).sum().item()
52         val_total += y.size(0)
53
54     # Compute metrics
55     train_loss /= len(train_loader)
56     train_acc = train_correct / train_total
57     val_loss /= len(val_loader)
58     val_acc = val_correct / val_total
59
60     # Log to TensorBoard
61     writer.add_scalar('Loss/train', train_loss, epoch)
62     writer.add_scalar('Loss/val', val_loss, epoch)
63     writer.add_scalar('Accuracy/train', train_acc, epoch)
64     writer.add_scalar('Accuracy/val', val_acc, epoch)
65     writer.add_scalar('LR', optimizer.param_groups[0]['lr'],
66 epoch)
67     writer.add_scalar('GradNorm', grad_norm, epoch)
68
69     # Print
70     print(f"Epoch {epoch}: Train Loss={train_loss:.4f}, "
71           f"Val Loss={val_loss:.4f}, Val Acc={val_acc:.4f}")
72
73     # Step scheduler
74     scheduler.step()
75 writer.close()

```

13.4.2 Recognizing Training Issues from Curves

Healthy training:

Train loss: Decreasing smoothly
 Val loss: Decreasing, tracking train loss
 Gap: Small (< 10% difference)

Overfitting:

Train loss: Very low
 Val loss: High and increasing
 Gap: Large and growing
 → Solution: More regularization, more data

Underfitting:

Train loss: High
 Val loss: High (similar to train)
 Gap: Small
 → Solution: Bigger model, train longer

Unstable training:

Both losses: Oscillating wildly
 → Solution: Lower learning rate, gradient clipping

13.5 Model Selection and Evaluation

13.5.1 Proper Train/Val/Test Split

```

1  # WRONG: Tune on test set
2  # This gives overly optimistic results!
3
4  # RIGHT: Three separate sets
5  # 1. Training set (70%): Train model
6  # 2. Validation set (15%): Tune hyperparameters
7  # 3. Test set (15%): Final evaluation (ONCE!)
8
9  from sklearn.model_selection import train_test_split
10
11 # Split data
12 X_train, X_temp, y_train, y_temp = train_test_split(
13     X, y, test_size=0.3, stratify=y, random_state=42
14 )
15 X_val, X_test, y_val, y_test = train_test_split(
16     X_temp, y_temp, test_size=0.5, stratify=y_temp,
17     random_state=42
18 )
19 # Training process:
20 # 1. Train many models with different hyperparameters
21 # 2. Evaluate all on validation set
22 # 3. Pick best model
23 # 4. Train best model on train+val (optional)
24 # 5. Evaluate ONCE on test set
25 # 6. Report test performance

```

13.5.2 Cross-Validation for Small Datasets

```

1  from sklearn.model_selection import KFold
2
3  kfold = KFold(n_splits=5, shuffle=True, random_state=42)
4
5  scores = []
6
7  for fold, (train_idx, val_idx) in enumerate(kfold.split(X)):
8      print(f"Fold {fold + 1}")
9
10     # Split data
11     X_train, X_val = X[train_idx], X[val_idx]
12     y_train, y_val = y[train_idx], y[val_idx]
13
14     # Create datasets
15     train_dataset = YourDataset(X_train, y_train)
16     val_dataset = YourDataset(X_val, y_val)
17
18     # Train model (reset for each fold!)
19     model = YourModel()
20     optimizer = torch.optim.Adam(model.parameters())
21
22     # Train...
23     val_score = train_and_evaluate(model, train_dataset,
24     val_dataset)
25     scores.append(val_score)
26
27 # Report mean and std
28 print(f"CV Score: {np.mean(scores):.4f}    {np.std(scores):.4f}")

```


13.6 Practical Tips

13.6.1 Start with Good Defaults

```

1  # Proven configuration for most tasks:
2
3  optimizer = torch.optim.AdamW(
4      model.parameters(),
5      lr=1e-3,           # Will likely need tuning
6      weight_decay=1e-4,
7      betas=(0.9, 0.999)
8  )
9
10 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(
11     optimizer, T_max=num_epochs
12 )
13
14 # Initialization: PyTorch defaults are usually good
15 # But for RNNs, use orthogonal init for hidden-to-hidden
16
17 # Batch size: 64 or 128
18 # Epochs: Start with 100, use early stopping
19 # Gradient clipping: max_norm=1.0 (for RNNs)

```

13.6.2 Hyperparameter Search Strategies

Grid search (exhaustive):

```

1  lrs = [1e-4, 1e-3, 1e-2]
2  weight_decays = [0, 1e-5, 1e-4, 1e-3]
3  dropouts = [0.2, 0.5]
4
5  best_val_loss = float('inf')
6  best_config = None
7
8  for lr in lrs:
9      for wd in weight_decays:
10         for dropout in dropouts:
11             val_loss = train_and_evaluate(lr, wd, dropout)
12
13             if val_loss < best_val_loss:
14                 best_val_loss = val_loss
15                 best_config = (lr, wd, dropout)
16
17  print(f"Best: LR={best_config[0]}, WD={best_config[1]}, "
18        f"Dropout={best_config[2]}")

```

Random search (more efficient):

```

1  import random
2
3  def random_config():
4      lr = 10 ** random.uniform(-5, -2)  # 1e-5 to 1e-2
5      wd = 10 ** random.uniform(-6, -3)  # 1e-6 to 1e-3
6      dropout = random.uniform(0.1, 0.7)
7      return lr, wd, dropout
8
9  num_trials = 20
10 results = []
11
12 for trial in range(num_trials):

```



```

13     lr, wd, dropout = random_config()
14     val_loss = train_and_evaluate(lr, wd, dropout)
15     results.append((val_loss, lr, wd, dropout))
16
17     # Sort by validation loss
18     results.sort()
19     best = results[0]
20     print(f"Best: Val Loss={best[0]:.4f}, LR={best[1]:.2e}, "
21           f"WD={best[2]:.2e}, Dropout={best[3]:.2f}")

```

Bayesian optimization (most efficient):

```

1  # Use libraries like Optuna or Ray Tune
2  import optuna
3
4  def objective(trial):
5      # Suggest hyperparameters
6      lr = trial.suggest_loguniform('lr', 1e-5, 1e-2)
7      wd = trial.suggest_loguniform('weight_decay', 1e-6, 1e-3)
8      dropout = trial.suggest_uniform('dropout', 0.1, 0.7)
9
10     # Train and return validation loss
11     val_loss = train_and_evaluate(lr, wd, dropout)
12     return val_loss
13
14     # Optimize
15     study = optuna.create_study(direction='minimize')
16     study.optimize(objective, n_trials=50)
17
18     print(f"Best hyperparameters: {study.best_params}")
19     print(f"Best validation loss: {study.best_value:.4f}")

```

13.7 Key Takeaways

Hyperparameter priorities:

1. Learning rate (MOST IMPORTANT)
2. Batch size
3. Architecture (depth, width)
4. Regularization (dropout, weight decay)
5. Everything else

Learning rate strategies:

- Use LR finder to find good starting point
- Start high, decay over time
- Cosine annealing works well
- Warmup helps for Transformers
- ReduceLROnPlateau is safe fallback

Batch size selection:

- Start with 32-128
- Larger batch = faster but may hurt generalization
- Scale LR with batch size (approximately)
- Use gradient accumulation if memory limited

Regularization:

- Dropout: 0.5 for FC layers, 0.1-0.2 for conv layers
- Weight decay: $1e-5$ to $1e-3$
- Data augmentation: Always for images
- Early stopping: Always use

Training monitoring:

- Track train/val loss and metrics
- Use TensorBoard for visualization
- Watch for overfitting (gap between train/val)
- Watch for underfitting (both losses high)

Model selection:

- Always use separate test set
- Tune on validation set only
- Test set evaluated ONCE at the end
- Use cross-validation for small datasets

Hyperparameter search:

- Random search better than grid search
- Bayesian optimization most efficient
- Start with good defaults
- Log everything

Common mistakes:

- Not tuning learning rate
- Evaluating on test set during development
- Using validation statistics to normalize test data
- Not using early stopping
- Training for too many epochs

14 Performance & Optimization

14.1 Introduction: Making Training Faster

Training can be slow. A model that takes 1 hour vs 10 hours per epoch makes a huge difference in development speed.

This section covers:

- Mixed precision training (2× speedup)
- GPU optimization
- DataLoader efficiency
- Memory management
- Profiling tools

14.2 Mixed Precision Training

Idea: Use FP16 (half precision) instead of FP32 (full precision) for most operations.

Benefits:

- 2-3× faster training
- 2× less memory (can use larger batch sizes)
- Minimal accuracy loss

14.2.1 Automatic Mixed Precision (AMP)

PyTorch provides automatic mixed precision via `torch.cuda.amp`:

```
1 from torch.cuda.amp import autocast, GradScaler
2
3 model = YourModel().cuda()
4 optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
5
6 # Create gradient scaler
7 scaler = GradScaler()
8
9 for epoch in range(num_epochs):
10     for x, y in train_loader:
11         x, y = x.cuda(), y.cuda()
12
13         optimizer.zero_grad()
14
15         # Forward pass in mixed precision
16         with autocast():
17             output = model(x)
18             loss = criterion(output, y)
19
20         # Backward pass with gradient scaling
21         scaler.scale(loss).backward()
22
23         # Unscale gradients and step
24         scaler.step(optimizer)
25         scaler.update()
```

How it works:

1. `autocast()`: Automatically chooses FP16 or FP32 for each operation
2. `GradScaler`: Scales loss to prevent underflow in FP16 gradients

[When to Use Mixed Precision] **Use AMP when:**

- Training on modern GPUs (Volta, Turing, Ampere)
- You want faster training
- You're memory-limited

Don't use if:

- GPU doesn't support FP16 (e.g., older GPUs)
- Model has numerical instability issues
- You need exact FP32 precision (rare)

Typical speedup: 1.5-3× depending on model and GPU.

14.2.2 Gradient Checkpointing

Trade computation for memory:

```

1 from torch.utils.checkpoint import checkpoint
2
3 class MemoryEfficientBlock(nn.Module):
4     def __init__(self, ...):
5         super().__init__()
6         self.layer = ExpensiveLayer()
7
8     def forward(self, x):
9         # Use checkpointing for this layer
10        return checkpoint(self.layer, x)
11
12 # Saves activations strategically
13 # Recomputes forward pass during backward (uses more time,
    less memory)

```

When to use:

- Very deep models (100+ layers)
- Running out of GPU memory
- Willing to trade 20-30% slower training for 2× memory reduction

14.3 GPU Optimization

14.3.1 Data Transfer Optimization

Pin memory for faster CPU→GPU transfer:

```

1 train_loader = DataLoader(
2     dataset,
3     batch_size=128,
4     num_workers=4,
5     pin_memory=True # Faster transfers
6 )
7
8 # Also call .cuda(non_blocking=True)
9 for x, y in train_loader:
10     x = x.cuda(non_blocking=True)
11     y = y.cuda(non_blocking=True)
12
13     # Model forward...
```

Keep data on GPU when possible:

```

1 # BAD: Moving data back and forth
2 for x, y in train_loader:
3     x = x.cuda()
4     output = model(x)
5     output = output.cpu() # Unnecessary transfer!
6     loss = criterion(output, y.cpu()) # Another transfer!
7
8 # GOOD: Keep everything on GPU
9 for x, y in train_loader:
10     x, y = x.cuda(), y.cuda()
11     output = model(x) # Stays on GPU
12     loss = criterion(output, y) # All on GPU
```

14.3.2 Batch Size and GPU Utilization

```

1 import torch
2
3 # Check GPU utilization
4 # nvidia-smi in terminal
5
6 # If GPU utilization < 80%, try:
7 # 1. Increase batch size
8 train_loader = DataLoader(dataset, batch_size=256) # Larger
9
10 # 2. Increase num_workers
11 train_loader = DataLoader(dataset, batch_size=128,
12                             num_workers=8)
13
14 # 3. Use larger model (if not memory-limited)
```

14.3.3 Multi-GPU Training

DataParallel (simple but not recommended):

```

1 if torch.cuda.device_count() > 1:
2     print(f"Using {torch.cuda.device_count()} GPUs")
3     model = nn.DataParallel(model)
4
```

```

5 model = model.cuda()
6
7 # Training as usual
8 # Batch automatically split across GPUs

```

DistributedDataParallel (better, recommended):

```

1 import torch.distributed as dist
2 from torch.nn.parallel import DistributedDataParallel as DDP
3
4 def setup(rank, world_size):
5     dist.init_process_group("nccl", rank=rank, world_size=
6         world_size)
7
8 def cleanup():
9     dist.destroy_process_group()
10
11 def train(rank, world_size):
12     setup(rank, world_size)
13
14     # Create model and move to GPU
15     model = YourModel().to(rank)
16     model = DDP(model, device_ids=[rank])
17
18     # Create distributed sampler
19     sampler = torch.utils.data.distributed.DistributedSampler(
20         dataset, num_replicas=world_size, rank=rank
21     )
22
23     train_loader = DataLoader(dataset, batch_size=128,
24         sampler=sampler)
25
26     # Training loop as usual
27     for epoch in range(num_epochs):
28         sampler.set_epoch(epoch) # Important for shuffling!
29         for x, y in train_loader:
30             # Training step...
31             pass
32
33     cleanup()
34
35 # Launch with torch.multiprocessing
36 import torch.multiprocessing as mp
37 mp.spawn(train, args=(world_size,), nprocs=world_size, join=
38     True)
39
40 # Or use torchrun:
41 # torchrun --nproc_per_node=4 train.py

```

DDP vs DataParallel:

- DDP: Faster, more efficient, recommended
- DP: Simpler but slower, single-process

14.4 DataLoader Optimization

14.4.1 Optimal Number of Workers

```

1 import time
2
3 def benchmark_dataloader(num_workers):
4     """Find optimal num_workers."""
5     loader = DataLoader(dataset, batch_size=128,
6                         num_workers=num_workers, pin_memory=
7                         True)
8     start = time.time()
9     for _ in loader:
10         pass # Just iterate
11     elapsed = time.time() - start
12
13     return elapsed
14
15 # Test different values
16 for num_workers in [0, 2, 4, 8, 16]:
17     elapsed = benchmark_dataloader(num_workers)
18     print(f"num_workers={num_workers}: {elapsed:.2f}s")
19
20 # Typical optimal: 4-8 workers
21 # More workers = more CPU usage but faster
22 # Too many workers = overhead dominates

```

Guidelines:

- Start with num_workers=4
- If CPU usage low: increase
- If CPU usage 100%: decrease
- Typical range: 2-8 workers

14.4.2 Prefetching

PyTorch DataLoader automatically prefetches, but you can optimize:

```

1 # Use persistent_workers to avoid recreating workers
2 train_loader = DataLoader(
3     dataset,
4     batch_size=128,
5     num_workers=4,
6     pin_memory=True,
7     persistent_workers=True, # Keep workers alive between
8     epochs
9     prefetch_factor=2 # Number of batches to prefetch per
10    worker
11 )

```

14.5 Memory Management

14.5.1 Monitoring Memory Usage

```

1 import torch
2
3 def print_memory_usage():
4     """Print current GPU memory usage."""
5     if torch.cuda.is_available():

```

```

6         allocated = torch.cuda.memory_allocated() / 1e9
7         reserved = torch.cuda.memory_reserved() / 1e9
8         print(f"Allocated: {allocated:.2f} GB")
9         print(f"Reserved: {reserved:.2f} GB")
10
11     # Check memory before and after operations
12     print("Before model creation:")
13     print_memory_usage()
14
15     model = LargeModel().cuda()
16
17     print("After model creation:")
18     print_memory_usage()

```

14.5.2 Reducing Memory Usage

1. Use gradient accumulation instead of large batches:

```

1  # Instead of batch_size=512 (OOM)
2  # Use batch_size=128 with accumulation
3
4  accumulation_steps = 4
5  effective_batch_size = 128 * 4  # = 512
6
7  optimizer.zero_grad()
8
9  for i, (x, y) in enumerate(train_loader):
10     output = model(x)
11     loss = criterion(output, y)
12     loss = loss / accumulation_steps  # Scale loss
13     loss.backward()  # Accumulate gradients
14
15     if (i + 1) % accumulation_steps == 0:
16         optimizer.step()
17         optimizer.zero_grad()

```

2. Clear cache periodically:

```

1  # If encountering OOM during long training
2  if epoch % 10 == 0:
3     torch.cuda.empty_cache()

```

3. Delete unnecessary variables:

```

1  # In evaluation, don't keep gradients
2  with torch.no_grad():
3     for x, y in val_loader:
4         output = model(x)
5         # ...
6
7  # If you need to keep large tensors temporarily
8  large_tensor = compute_something()
9  result = process(large_tensor)
10 del large_tensor  # Free memory immediately
11 torch.cuda.empty_cache()

```

4. Use smaller data types:

```

1  # FP32 (default)
2  x = torch.randn(1000, 1000)  # 4 MB

```



```
3
4 # FP16 (half)
5 x = torch.randn(1000, 1000, dtype=torch.float16) # 2 MB
6
7 # For labels, use appropriate type
8 labels = torch.randint(0, 10, (1000,)), dtype=torch.long) #
    Not float!
```

14.6 Profiling and Bottleneck Analysis

14.6.1 PyTorch Profiler

```

1 from torch.profiler import profile, ProfilerActivity
2
3 model = YourModel().cuda()
4
5 # Profile training
6 with profile(
7     activities=[ProfilerActivity.CPU, ProfilerActivity.CUDA],
8     record_shapes=True,
9     profile_memory=True,
10    with_stack=True
11 ) as prof:
12     for i, (x, y) in enumerate(train_loader):
13         if i >= 10: # Profile first 10 batches
14             break
15
16         x, y = x.cuda(), y.cuda()
17
18         optimizer.zero_grad()
19         output = model(x)
20         loss = criterion(output, y)
21         loss.backward()
22         optimizer.step()
23
24 # Print profiling results
25 print(prof.key_averages().table(sort_by="cuda_time_total",
26                                row_limit=10))
27
28 # Export for Chrome tracing
29 prof.export_chrome_trace("trace.json")
30 # View at chrome://tracing

```

14.6.2 Simple Timing

```

1 import time
2
3 # Time entire epoch
4 start = time.time()
5 for x, y in train_loader:
6     # Training step...
7     pass
8 elapsed = time.time() - start
9 print(f"Epoch time: {elapsed:.2f}s")
10
11 # Time individual operations
12 torch.cuda.synchronize() # Wait for GPU
13 start = time.time()
14
15 output = model(x)
16
17 torch.cuda.synchronize()
18 elapsed = time.time() - start
19 print(f"Forward pass: {elapsed*1000:.2f}ms")

```

14.6.3 Identify Bottlenecks

```

1 def profile_training_loop(model, train_loader, num_batches
2                             =10):

```

```

2  """Profile different parts of training."""
3
4  times = {
5      'data_loading': 0,
6      'to_gpu': 0,
7      'forward': 0,
8      'backward': 0,
9      'optimizer': 0,
10     'total': 0
11 }
12
13 total_start = time.time()
14 data_start = time.time()
15
16 for i, (x, y) in enumerate(train_loader):
17     if i >= num_batches:
18         break
19
20     times['data_loading'] += time.time() - data_start
21
22     # Transfer to GPU
23     gpu_start = time.time()
24     x, y = x.cuda(), y.cuda()
25     torch.cuda.synchronize()
26     times['to_gpu'] += time.time() - gpu_start
27
28     # Forward
29     forward_start = time.time()
30     output = model(x)
31     torch.cuda.synchronize()
32     times['forward'] += time.time() - forward_start
33
34     # Backward
35     backward_start = time.time()
36     loss = criterion(output, y)
37     loss.backward()
38     torch.cuda.synchronize()
39     times['backward'] += time.time() - backward_start
40
41     # Optimizer
42     opt_start = time.time()
43     optimizer.step()
44     optimizer.zero_grad()
45     torch.cuda.synchronize()
46     times['optimizer'] += time.time() - opt_start
47
48     data_start = time.time()
49
50 times['total'] = time.time() - total_start
51
52 # Print breakdown
53 print("Time breakdown:")
54 for key, value in times.items():
55     if key != 'total':
56         pct = (value / times['total']) * 100
57         print(f"  {key}: {value:.3f}s ({pct:.1f}%)"
58     else:
59         print(f"  {key}: {value:.3f}s")
60
61 # Usage
62 profile_training_loop(model, train_loader)
63 """

```

```
64 Time breakdown:
65   data_loading: 0.421s (25.3%)
66   to_gpu: 0.089s (5.3%)
67   forward: 0.567s (34.1%)
68   backward: 0.423s (25.4%)
69   optimizer: 0.165s (9.9%)
70   total: 1.665s
71   """
72
73   # If data_loading is high: increase num_workers
74   # If forward/backward is high: model is slow (expected)
75   # If to_gpu is high: use pin_memory, non_blocking transfers
```

14.7 Production Deployment

14.7.1 Model Export

TorchScript (JIT compilation):

```

1 # Option 1: Tracing
2 model = YourModel()
3 model.eval()
4
5 example_input = torch.randn(1, 3, 224, 224)
6 traced_model = torch.jit.trace(model, example_input)
7
8 # Save
9 traced_model.save("model_traced.pt")
10
11 # Load and use
12 loaded_model = torch.jit.load("model_traced.pt")
13 output = loaded_model(example_input)

```

```

1 # Option 2: Scripting (handles control flow)
2 scripted_model = torch.jit.script(model)
3 scripted_model.save("model_scripted.pt")

```

ONNX (cross-framework):

```

1 import torch.onnx
2
3 model = YourModel()
4 model.eval()
5
6 dummy_input = torch.randn(1, 3, 224, 224)
7
8 torch.onnx.export(
9     model,
10    dummy_input,
11    "model.onnx",
12    export_params=True,
13    opset_version=11,
14    input_names=['input'],
15    output_names=['output'],
16    dynamic_axes={'input': {0: 'batch_size'},
17                  'output': {0: 'batch_size'}}
18 )
19
20 # Can now use in other frameworks (TensorFlow, etc.)

```

14.7.2 Inference Optimization

```

1 # 1. Set to eval mode
2 model.eval()
3
4 # 2. Disable gradient computation
5 with torch.no_grad():
6     output = model(input)
7
8 # 3. Use FP16 if possible
9 with torch.cuda.amp.autocast():
10    output = model(input)
11
12 # 4. Batch inference

```

```

13 # Process multiple inputs together
14 batch = torch.stack([input1, input2, input3])
15 outputs = model(batch)
16
17 # 5. Use torch.jit for production
18 model = torch.jit.load("model.pt")

```

14.8 Key Takeaways

Mixed precision training:

- Use `torch.cuda.amp` for $2\times$ speedup
- Minimal code changes required
- Works on modern GPUs (Volta+)
- Typical speedup: $1.5\text{-}3\times$

GPU optimization:

- Keep data on GPU (avoid CPU-GPU transfers)
- Use `pin_memory=True` and `non_blocking=True`
- Maximize batch size (within memory limits)
- Use `DistributedDataParallel` for multi-GPU

DataLoader optimization:

- Start with `num_workers=4`
- Use `persistent_workers=True`
- Benchmark different worker counts
- Typical optimal: 4-8 workers

Memory management:

- Use gradient accumulation for large effective batch sizes
- Clear cache periodically if needed
- Delete unnecessary variables
- Use gradient checkpointing for very deep models

Profiling:

- Use PyTorch Profiler for detailed analysis
- Identify bottlenecks (data loading, forward, backward)
- Focus optimization efforts on slowest parts
- Profile regularly during development

Production deployment:

- Use TorchScript for inference
- Export to ONNX for cross-framework compatibility
- Always use `model.eval()` and `torch.no_grad()`
- Batch inference when possible
- Consider FP16 inference

Common performance issues:

- Data loading too slow \rightarrow increase `num_workers`

- GPU underutilized → increase batch size or model size
- Out of memory → reduce batch size, use gradient accumulation, or AMP
- Training too slow → use mixed precision, profile bottlenecks

Part IV

Integrative Challenges

15 Grand Challenges

15.1 Introduction: Putting It All Together

You've learned the fundamentals. Now it's time to combine everything into complete, real-world projects.

These challenges integrate multiple concepts:

- Multiple architecture types (CNNs, RNNs, Transformers)
- Proper training pipelines (data loading, optimization, monitoring)
- Best practices (regularization, debugging, performance)
- Scientific computing applications

How to approach these challenges:

1. Start simple (baseline model)
2. Get it working (overfit small data)
3. Make it good (full dataset + regularization)
4. Make it great (hyperparameter tuning + tricks)
5. Document everything (what worked, what didn't)

15.2 Challenge 1: Time Series Forecasting System

Difficulty: ★★★

Goal: Build a complete forecasting system for multivariate time series.

Dataset: Use real data (weather, stock prices, sensor readings) or generate complex synthetic data.

Requirements:

1. Data Pipeline

- Load and preprocess time series data
- Create sliding windows (past 100 steps → predict next 10)
- Handle missing values
- Normalize using training statistics
- Split into train/val/test (temporal split, not random!)

2. Model Architecture

Implement and compare three approaches:

- **LSTM-based:** Stacked LSTM with attention
- **CNN-based:** 1D CNN with residual connections
- **Transformer-based:** Encoder-only Transformer

3. Training

- Proper learning rate schedule (warmup + decay)

- Early stopping on validation set
- Gradient clipping (essential for RNNs)
- Mixed precision training
- TensorBoard logging

4. Evaluation

- Multiple metrics: MSE, MAE, MAPE
- One-step ahead vs multi-step ahead prediction
- Visualization: predicted vs actual
- Confidence intervals (optional)

5. Analysis

- Which architecture works best?
- How far ahead can you predict accurately?
- What patterns does the model learn?
- Where does it fail?

Starter template:

```

1 class TimeSeriesDataset(Dataset):
2     def __init__(self, data, window_size, forecast_horizon):
3         # Your implementation
4         pass
5
6     def __len__(self):
7         return len(self.data) - self.window_size - self.
forecast_horizon + 1
8
9     def __getitem__(self, idx):
10        # Return (input_window, target_window)
11        pass
12
13 class LSTMForecaster(nn.Module):
14     def __init__(self, input_size, hidden_size, num_layers,
forecast_horizon):
15         super().__init__()
16         self.lstm = nn.LSTM(input_size, hidden_size,
num_layers,
17                               batch_first=True)
18         self.fc = nn.Linear(hidden_size, forecast_horizon)
19
20     def forward(self, x):
21         # LSTM encoding
22         output, (h_n, c_n) = self.lstm(x)
23         # Predict from final hidden state
24         prediction = self.fc(h_n[-1])
25         return prediction
26
27 # Your training loop with all best practices

```

Bonus challenges:

- Add attention mechanism to LSTM
- Implement probabilistic forecasting (predict distribution, not point estimate)
- Handle multiple related time series (multivariate)
- Deploy as REST API

15.3 Challenge 2: Hybrid Vision-Language Model

Difficulty: ★★★★★

Goal: Build a model that combines visual and textual information.

Task: Image captioning or visual question answering.

Requirements:

1. Architecture

Implement encoder-decoder with:

- **Image encoder:** ResNet or Vision Transformer
- **Text decoder:** LSTM or Transformer decoder
- **Cross-modal attention:** Decoder attends to image features

```

1 class ImageCaptioningModel(nn.Module):
2     def __init__(self, vocab_size, embed_dim=512, hidden_dim
      =512):
3         super().__init__()
4
5         # Image encoder (pretrained ResNet)
6         resnet = models.resnet50(pretrained=True)
7         # Remove final classification layer
8         self.encoder = nn.Sequential(*list(resnet.children())
      [:-2])
9
10        # Freeze encoder (transfer learning)
11        for param in self.encoder.parameters():
12            param.requires_grad = False
13
14        # Project image features
15        self.image_proj = nn.Linear(2048, embed_dim)
16
17        # Text decoder
18        self.embedding = nn.Embedding(vocab_size, embed_dim)
19        self.lstm = nn.LSTM(embed_dim, hidden_dim, num_layers
      =2,
20                               batch_first=True)
21
22        # Attention mechanism
23        self.attention = nn.MultiheadAttention(embed_dim,
      num_heads=8)
24
25        # Output projection
26        self.fc_out = nn.Linear(hidden_dim, vocab_size)
27
28    def forward(self, images, captions):
29        # Encode images: (batch, 2048, 7, 7)
30        img_features = self.encoder(images)
31
32        # Reshape: (batch, 49, 2048)
33        batch_size = img_features.size(0)
34        img_features = img_features.view(batch_size, 2048,
      -1)
35        img_features = img_features.transpose(1, 2)
36
37        # Project: (batch, 49, embed_dim)
38        img_features = self.image_proj(img_features)
39
40        # Embed captions

```

```

41     caption_embeds = self.embedding(captions)
42
43     # LSTM decoding with attention
44     lstm_out, _ = self.lstm(caption_embeds)
45
46     # Attend to image features
47     attn_out, _ = self.attention(lstm_out, img_features,
img_features)
48
49     # Combine and project
50     combined = lstm_out + attn_out
51     output = self.fc_out(combined)
52
53     return output

```

2. Data Preparation

- Use COCO dataset or similar
- Build vocabulary from captions
- Tokenize captions
- Implement proper data augmentation for images
- Handle variable-length captions (padding + masking)

3. Training

- Teacher forcing during training
- Proper cross-entropy loss (ignore padding tokens)
- Learning rate warmup + decay
- Mixed precision training
- Gradient accumulation (if memory limited)

4. Generation

Implement multiple decoding strategies:

```

1  def generate_caption(model, image, vocab, max_len=20,
2                        strategy='greedy', temperature=1.0,
3                        beam_width=5):
4      """
5      Generate caption from image.
6
7      Args:
8          strategy: 'greedy', 'sampling', or 'beam_search'
9      """
10     model.eval()
11
12     # Encode image
13     with torch.no_grad():
14         img_features = model.encode_image(image)
15
16     # Start token
17     caption = [vocab['<start>']]
18
19     if strategy == 'greedy':
20         # Greedy decoding (take argmax at each step)
21         for _ in range(max_len):
22             logits = model.decode_step(img_features, caption)
23             next_word = logits.argmax()
24             caption.append(next_word.item())

```

```

24         if next_word == vocab['<end>']:
25             break
26
27     elif strategy == 'sampling':
28         # Sample from distribution
29         for _ in range(max_len):
30             logits = model.decode_step(img_features, caption)
31             logits = logits / temperature
32             probs = F.softmax(logits, dim=-1)
33             next_word = torch.multinomial(probs, 1)
34             caption.append(next_word.item())
35             if next_word == vocab['<end>']:
36                 break
37
38     elif strategy == 'beam_search':
39         # Beam search (keep top-k candidates)
40         # Your implementation
41         pass
42
43     return caption

```

5. Evaluation

- BLEU score (text similarity metric)
- Qualitative: Visual inspection of generated captions
- Error analysis: Where does model fail?
- Attention visualization: What does model look at?

Bonus challenges:

- Implement beam search decoding
- Add image-text contrastive learning (CLIP-style)
- Fine-tune encoder (not just freeze)
- Add copy mechanism for rare words

15.4 Challenge 3: Denoising Scientific Data

Difficulty: ★★★

Goal: Build a U-Net style architecture to denoise 2D scientific data.

Application: Medical images, microscopy, sensor data, astronomical images.

Requirements:

1. Data Generation

Create noisy data:

```

1 def generate_noisy_data(clean_images, noise_level=0.1):
2     """Add Gaussian noise to clean images."""
3     noise = torch.randn_like(clean_images) * noise_level
4     noisy = clean_images + noise
5     return noisy, clean_images
6
7 # Use real dataset (medical images, etc.) or synthetic

```

2. U-Net Architecture

Implement complete U-Net with:

- Encoder: Downsampling path with residual blocks
- Decoder: Upsampling path with skip connections
- Attention gates (optional but recommended)
- Batch normalization throughout

```

1 class UNet(nn.Module):
2     """U-Net for image denoising/reconstruction."""
3
4     def __init__(self, in_channels=1, out_channels=1):
5         super().__init__()
6
7         # Encoder
8         self.enc1 = self.conv_block(in_channels, 64)
9         self.enc2 = self.conv_block(64, 128)
10        self.enc3 = self.conv_block(128, 256)
11        self.enc4 = self.conv_block(256, 512)
12
13        # Bottleneck
14        self.bottleneck = self.conv_block(512, 1024)
15
16        # Decoder with skip connections
17        self.upconv4 = nn.ConvTranspose2d(1024, 512, 2,
stride=2)
18        self.dec4 = self.conv_block(1024, 512) # 1024 = 512
+ 512 from skip
19
20        self.upconv3 = nn.ConvTranspose2d(512, 256, 2, stride
=2)
21        self.dec3 = self.conv_block(512, 256)
22
23        self.upconv2 = nn.ConvTranspose2d(256, 128, 2, stride
=2)
24        self.dec2 = self.conv_block(256, 128)
25
26        self.upconv1 = nn.ConvTranspose2d(128, 64, 2, stride
=2)

```

```

27     self.dec1 = self.conv_block(128, 64)
28
29     # Output
30     self.out = nn.Conv2d(64, out_channels, 1)
31
32     self.pool = nn.MaxPool2d(2, 2)
33
34     def conv_block(self, in_ch, out_ch):
35         """Double convolution block."""
36         return nn.Sequential(
37             nn.Conv2d(in_ch, out_ch, 3, padding=1),
38             nn.BatchNorm2d(out_ch),
39             nn.ReLU(inplace=True),
40             nn.Conv2d(out_ch, out_ch, 3, padding=1),
41             nn.BatchNorm2d(out_ch),
42             nn.ReLU(inplace=True)
43         )
44
45     def forward(self, x):
46         # Encoder
47         enc1 = self.enc1(x)
48         enc2 = self.enc2(self.pool(enc1))
49         enc3 = self.enc3(self.pool(enc2))
50         enc4 = self.enc4(self.pool(enc3))
51
52         # Bottleneck
53         bottleneck = self.bottleneck(self.pool(enc4))
54
55         # Decoder with skip connections
56         dec4 = self.upconv4(bottleneck)
57         dec4 = torch.cat([dec4, enc4], dim=1)
58         dec4 = self.dec4(dec4)
59
60         dec3 = self.upconv3(dec4)
61         dec3 = torch.cat([dec3, enc3], dim=1)
62         dec3 = self.dec3(dec3)
63
64         dec2 = self.upconv2(dec3)
65         dec2 = torch.cat([dec2, enc2], dim=1)
66         dec2 = self.dec2(dec2)
67
68         dec1 = self.upconv1(dec2)
69         dec1 = torch.cat([dec1, enc1], dim=1)
70         dec1 = self.dec1(dec1)
71
72         return self.out(dec1)

```

3. Training

- Loss function: MSE + Perceptual loss (optional)
- Data augmentation: Random flips, rotations
- Learning rate schedule
- Monitor both MSE and PSNR metrics

4. Evaluation

```

1 def calculate_psnr(img1, img2):
2     """Calculate Peak Signal-to-Noise Ratio."""
3     mse = torch.mean((img1 - img2) ** 2)
4     if mse == 0:
5         return float('inf')

```

```

6     max_pixel = 1.0
7     psnr = 20 * torch.log10(max_pixel / torch.sqrt(mse))
8     return psnr.item()
9
10    def calculate_ssim(img1, img2):
11        """Structural Similarity Index (use skimage or pytorch).
12        """
13        from skimage.metrics import structural_similarity
14        return structural_similarity(
15            img1.cpu().numpy(), img2.cpu().numpy(),
16            data_range=1.0
17        )
18
19    # Evaluate
20    model.eval()
21    psnrs = []
22    ssims = []
23
24    with torch.no_grad():
25        for noisy, clean in test_loader:
26            denoised = model(noisy)
27
28            for i in range(len(noisy)):
29                psnr = calculate_psnr(denoised[i], clean[i])
30                ssim = calculate_ssim(denoised[i], clean[i])
31                psnrs.append(psnr)
32                ssims.append(ssim)
33
34    print(f"Average PSNR: {np.mean(psnrs):.2f} dB")
35    print(f"Average SSIM: {np.mean(ssims):.4f}")

```

5. Visualization

```

1    def visualize_denoising(model, noisy_images, clean_images):
2        """Visualize denoising results."""
3        model.eval()
4
5        with torch.no_grad():
6            denoised = model(noisy_images)
7
8        fig, axes = plt.subplots(len(noisy_images), 3, figsize
9                                =(12, 4*len(noisy_images)))
10
11        for i in range(len(noisy_images)):
12            # Noisy
13            axes[i, 0].imshow(noisy_images[i].squeeze().cpu(),
14                             cmap='gray')
15            axes[i, 0].set_title('Noisy')
16            axes[i, 0].axis('off')
17
18            # Denoised
19            axes[i, 1].imshow(denoised[i].squeeze().cpu(), cmap='
20            gray')
21            axes[i, 1].set_title('Denoised')
22            axes[i, 1].axis('off')
23
24            # Clean
25            axes[i, 2].imshow(clean_images[i].squeeze().cpu(),
26                             cmap='gray')
27            axes[i, 2].set_title('Clean (Target)')
28            axes[i, 2].axis('off')

```

```
26 plt.tight_layout()  
27 plt.show()
```

Bonus challenges:

- Add attention gates between encoder and decoder
- Implement progressive training (start with low resolution)
- Handle different noise types (Gaussian, salt-and-pepper, Poisson)
- Add perceptual loss using pretrained VGG

15.5 Challenge 4: Custom Transformer for Sequence Tasks

Difficulty: ★★★★★

Goal: Build a complete Transformer from scratch for sequence classification or generation.

Task: Text classification, sequence generation, or time series classification.

Requirements:

1. Complete Transformer Implementation

Build all components from scratch (no `nn.Transformer`):

- Multi-head attention
- Position-wise feed-forward
- Positional encoding
- Layer normalization
- Encoder and/or decoder blocks

2. Training Infrastructure

- Warmup learning rate schedule
- Gradient clipping
- Label smoothing
- Mixed precision training
- Checkpointing and recovery

3. Advanced Features

Implement at least two of:

- Learning rate warmup + cosine annealing
- Custom learning rate schedule
- Gradient accumulation
- Model ensemble
- Knowledge distillation

4. Comprehensive Evaluation

- Validation metrics
- Attention visualization
- Error analysis
- Compare with baseline (LSTM, simple MLP)
- Ablation studies (remove components, measure impact)

5. Analysis

Answer these questions:

- How many heads are optimal?
- What does each attention head learn?
- How deep should the model be?
- What's the effect of positional encoding?
- Where does the model fail?

15.6 Challenge 5: End-to-End ML Pipeline

Difficulty: ★★★★★

Goal: Build a production-ready machine learning system.

Task: Choose your own (classification, regression, generation).

Requirements:

1. Complete Pipeline

```
project/
  data/
    raw/           # Raw data
    processed/     # Preprocessed data
    preprocessing.py # Data processing scripts
  models/
    architectures.py # Model definitions
    pretrained/     # Saved models
  training/
    train.py        # Training script
    evaluate.py     # Evaluation script
    utils.py        # Training utilities
  configs/
    base_config.yaml # Base configuration
    experiment_*.yaml # Experiment configs
  notebooks/
    exploration.ipynb # Data exploration
    analysis.ipynb   # Results analysis
  tests/
    test_*.py        # Unit tests
  requirements.txt
  README.md
  main.py            # Entry point
```

2. Code Quality

- Type hints
- Docstrings for all functions
- Unit tests for critical functions
- Configuration via YAML
- Logging (not just print statements)
- Error handling

3. Experiment Tracking

- TensorBoard or Weights & Biases
- Log hyperparameters
- Log metrics over time
- Save model checkpoints
- Version control (Git)

4. Reproducibility

- Set random seeds
- Save exact hyperparameters

- Save training logs
- Document environment (requirements.txt)
- Save data preprocessing steps

5. Deployment

- Export model (TorchScript or ONNX)
- Create inference script
- Docker container (optional)
- REST API using FastAPI (optional)
- Documentation for usage

Evaluation criteria:

- Model performance
- Code quality and organization
- Reproducibility
- Documentation
- Completeness of pipeline

15.7 Evaluation and Next Steps

How to know you've succeeded:

1. **It works:** Model trains without errors
2. **It learns:** Loss decreases, metrics improve
3. **It generalizes:** Validation performance is good
4. **You understand it:** Can explain why it works
5. **It's reproducible:** Others can replicate your results

Beyond this guide:

Keep learning:

- Read recent papers (arXiv, conferences)
- Implement papers from scratch
- Contribute to open source projects
- Join competitions (Kaggle, etc.)
- Build projects in your domain

Advanced topics to explore:

- Self-supervised learning
- Meta-learning
- Neural architecture search
- Efficient networks (quantization, pruning)
- Diffusion models
- Large language models
- Multimodal learning
- Reinforcement learning

Stay updated:

- Follow PyTorch blog and releases

- Join ML communities (Reddit, Discord, forums)
- Attend conferences (NeurIPS, ICML, ICLR)
- Read technical blogs
- Experiment with new techniques

15.8 Final Thoughts

You now have a comprehensive toolkit for deep learning with PyTorch. You've learned:

- **Foundations:** Tensors, autograd, modules, training loops, data loading
- **Architectures:** MLPs, CNNs, ResNets, RNNs, Transformers
- **Techniques:** Normalization, regularization, optimization
- **Practical skills:** Debugging, training best practices, performance optimization

Remember:

- Start simple, add complexity gradually
- Always check if your model can overfit small data
- Visualization and monitoring are essential
- Good engineering matters as much as good models
- Deep learning is iterative—experiment and learn

The best way to master these concepts is to **build things**. Pick a challenge that excites you and start coding!

Good luck, and happy building!