# PyTorch Functions: From Fundamentals to Neural Networks
### An Incremental Learning Guide

An Educational Tutorial

August 13, 2025

# Contents

# Chapter 1

# Mathematical Foundations

## 1.1 Linear Algebra with PyTorch

### 1.1.1 torch.linalg Functions

**Purpose:** PyTorch's linear algebra operations for mathematical computations in deep learning.

**Simple Example:**

**Tensor Shape Transformation Visualization:**

```
┌─────────────────────────┐
│       1D Tensor          │      Shape: (12,)
│ ,2,3,4,5,6,7,8,9,10,11   │
└─────────────────────────┘
            │
            ▼
    ┌─────────────────┐
    │    2D Tensor     │          Shape: (3,4)
    │   3×4 matrix     │
    └─────────────────┘
              │
              ▼
      ┌───────────────┐
      │      3D        │          Shape: (2,2,3)
      │    Tensor      │
      │    2×2×3       │
      └───────────────┘
```

```python
import torch
import torch.linalg as linalg

# Matrix operations
A = torch.randn(3, 3)
B = torch.randn(3, 3)

# Matrix multiplication
C = torch.matmul(A, B)  # or A @ B
print(f"Matrix product shape: {C.shape}")
# Output: Matrix product shape: torch.Size([3, 3])

# Determinant
det_A = torch.linalg.det(A)
```

```python
print(f"Determinant: {det_A}")
# Output: Determinant: tensor(-1.2354)

# Matrix inverse
A_inv = torch.linalg.inv(A)
print(f"Inverse verification: {torch.allclose(A @ A_inv, torch.eye(3))}")
# Output: Inverse verification: True

# Eigenvalues and eigenvectors
eigenvals, eigenvecs = torch.linalg.eig(A)
print(f"Eigenvalues: {eigenvals}")
# Output: Eigenvalues: tensor([-1.5678+0.0000j,  0.8901+1.2345j,  0.8901-1.2345j])
print(f"Eigenvectors shape: {eigenvecs.shape}")
# Output: Eigenvectors shape: torch.Size([3, 3])

# Singular Value Decomposition (SVD)
U, S, Vh = torch.linalg.svd(A)
print(f"SVD shapes: U{U.shape}, S{S.shape}, Vh{Vh.shape}")
# Output: SVD shapes: U torch.Size([3, 3]), S torch.Size([3]), Vh torch.Size([3,
#  3])

# Matrix norm
frobenius_norm = torch.linalg.norm(A, ord='fro')
nuclear_norm = torch.linalg.norm(A, ord='nuc')
print(f"Frobenius norm: {frobenius_norm}")
# Output: Frobenius norm: tensor(3.1623)
print(f"Nuclear norm: {nuclear_norm}")
# Output: Nuclear norm: tensor(4.5678)
```

**Complex Example - Principal Component Analysis:**

```python
def pca_torch(X, n_components):
    """
    Principal Component Analysis using PyTorch
    X: (n_samples, n_features)
    """
    # Center the data
    X_centered = X - X.mean(dim=0)

    # Compute covariance matrix
    n_samples = X.size(0)
    cov_matrix = (X_centered.T @ X_centered) / (n_samples - 1)

    # Eigendecomposition
    eigenvals, eigenvecs = torch.linalg.eigh(cov_matrix)

    # Sort eigenvalues and eigenvectors in descending order
```

```
17        idx = torch.argsort(eigenvals, descending=True)
18        eigenvals = eigenvals[idx]
19        eigenvecs = eigenvecs[:, idx]
20
21        # Select top n_components
22        components = eigenvecs[:, :n_components]
23        explained_variance = eigenvals[:n_components]
24
25        # Transform data
26        X_pca = X_centered @ components
27
28        return X_pca, components, explained_variance
29
30   # Example usage
31   data = torch.randn(100, 50)  # 100 samples, 50 features
32   X_reduced, components, var_explained = pca_torch(data, n_components=10)
33   print(f"Original shape: {data.shape}")
34   # Output: Original shape: torch.Size([100, 4])
35   print(f"Reduced shape: {X_reduced.shape}")
36   # Output: Reduced shape: torch.Size([100, 2])
37   print(f"Explained variance ratio: {var_explained / var_explained.sum()}")
38   # Output: Explained variance ratio: tensor([0.7296, 0.2277])
```

## 1.2 Probability Distributions

### 1.2.1 torch.distributions

**Purpose:** Probability distributions for probabilistic modeling and sampling.

**Simple Example:**

```
1    import torch.distributions as dist
2
3    # Normal distribution
4    normal = dist.Normal(loc=0.0, scale=1.0)
5    samples = normal.sample((1000,))
6    log_probs = normal.log_prob(samples)
7
8    print(f"Sample mean: {samples.mean():.3f}")
9    # Output: Sample mean: 0.012
10   print(f"Sample std: {samples.std():.3f}")
11   # Output: Sample std: 0.998
12
13   # Categorical distribution
14   categorical = dist.Categorical(probs=torch.tensor([0.1, 0.3, 0.6]))
15   cat_samples = categorical.sample((100,))
16   print(f"Categorical samples: {cat_samples[:10]}")
17   # Output: Categorical samples: tensor([2, 1, 2, 2, 0, 2, 1, 2, 2, 1])
```

```python
18
19   # Beta distribution
20   beta = dist.Beta(concentration1=2.0, concentration0=1.0)
21   beta_samples = beta.sample((100,))
22   print(f"Beta samples range: [{beta_samples.min():.3f}, {beta_samples.max():.3f}]")
23   # Output: Beta samples range: [0.126, 0.984]
24
25   # Multivariate Normal
26   mvn = dist.MultivariateNormal(
27       loc=torch.zeros(3),
28       covariance_matrix=torch.eye(3)
29   )
30   mvn_samples = mvn.sample((10,))
31   print(f"MVN samples shape: {mvn_samples.shape}")
32   # Output: MVN samples shape: torch.Size([10, 3])
```

**Complex Example - Variational Inference:**

```python
1    class VariationalBayesianLinear(nn.Module):
2        """Bayesian Linear Layer with Variational Inference"""
3        def __init__(self, in_features, out_features):
4            super().__init__()
5            self.in_features = in_features
6            self.out_features = out_features
7
8            # Weight parameters (mean and log variance)
9            self.weight_mu = nn.Parameter(torch.randn(out_features, in_features) *
             ↪   0.1)
10           self.weight_logvar = nn.Parameter(torch.randn(out_features, in_features) *
             ↪   0.1)
11
12           # Bias parameters
13           self.bias_mu = nn.Parameter(torch.randn(out_features) * 0.1)
14           self.bias_logvar = nn.Parameter(torch.randn(out_features) * 0.1)
15
16           # Prior distributions
17           self.weight_prior = dist.Normal(0, 1)
18           self.bias_prior = dist.Normal(0, 1)
19
20       def forward(self, x):
21           # Sample weights and biases
22           weight_std = torch.exp(0.5 * self.weight_logvar)
23           weight = dist.Normal(self.weight_mu, weight_std).rsample()
24
25           bias_std = torch.exp(0.5 * self.bias_logvar)
26           bias = dist.Normal(self.bias_mu, bias_std).rsample()
27
```

```python
28          return F.linear(x, weight, bias)

30      def kl_divergence(self):
31          """Compute KL divergence between posterior and prior"""
32          # Weight KL divergence
33          weight_posterior = dist.Normal(self.weight_mu, torch.exp(0.5 *
            ↪  self.weight_logvar))
34          weight_kl = dist.kl_divergence(weight_posterior, self.weight_prior).sum()

36          # Bias KL divergence
37          bias_posterior = dist.Normal(self.bias_mu, torch.exp(0.5 *
            ↪  self.bias_logvar))
38          bias_kl = dist.kl_divergence(bias_posterior, self.bias_prior).sum()

40          return weight_kl + bias_kl

42  # Example usage in a Bayesian Neural Network
43  class BayesianMLP(nn.Module):
44      def __init__(self, input_dim, hidden_dim, output_dim):
45          super().__init__()
46          self.layer1 = VariationalBayesianLinear(input_dim, hidden_dim)
47          self.layer2 = VariationalBayesianLinear(hidden_dim, output_dim)

49      def forward(self, x):
50          x = torch.relu(self.layer1(x))
51          return self.layer2(x)

53      def kl_divergence(self):
54          return self.layer1.kl_divergence() + self.layer2.kl_divergence()

56  # Training with ELBO (Evidence Lower BOund)
57  def train_bayesian_model(model, dataloader, epochs=10):
58      optimizer = torch.optim.Adam(model.parameters(), lr=0.01)

60      for epoch in range(epochs):
61          for batch_x, batch_y in dataloader:
62              optimizer.zero_grad()

64              # Forward pass
65              predictions = model(batch_x)

67              # Likelihood loss
68              likelihood_loss = F.mse_loss(predictions, batch_y)

70              # KL divergence
71              kl_loss = model.kl_divergence()

73              # ELBO = -likelihood + KL divergence
```

```
74              loss = likelihood_loss + kl_loss / len(dataloader.dataset)

75

76              loss.backward()
77              optimizer.step()

78

79          print(f"Epoch {epoch}: Loss = {loss.item():.4f}")
80          # Output: Epoch 0: Loss = 3.2456
```

## 1.3   Information Theory

### 1.3.1   Entropy and Mutual Information

**Purpose:** Information-theoretic measures for understanding learning and generalization.
      **Simple Example:**

```python
1   def entropy(probs, dim=-1):
2       """Compute entropy of probability distribution"""
3       # Add small epsilon to avoid log(0)
4       eps = 1e-8
5       return -torch.sum(probs * torch.log(probs + eps), dim=dim)

6

7   def cross_entropy(p, q, dim=-1):
8       """Cross entropy between distributions p and q"""
9       eps = 1e-8
10      return -torch.sum(p * torch.log(q + eps), dim=dim)

11

12  def kl_divergence(p, q, dim=-1):
13      """KL divergence between distributions p and q"""
14      return cross_entropy(p, q, dim) - entropy(p, dim)

15

16  # Example: Analyze model confidence
17  def analyze_model_uncertainty(model, dataloader):
18      model.eval()
19      entropies = []

20

21      with torch.no_grad():
22          for batch_x, _ in dataloader:
23              logits = model(batch_x)
24              probs = F.softmax(logits, dim=1)

25

26              # Compute entropy for each prediction
27              batch_entropy = entropy(probs, dim=1)
28              entropies.append(batch_entropy)

29

30      all_entropies = torch.cat(entropies)

31

32      print(f"Mean prediction entropy: {all_entropies.mean():.4f}")
```

```python
33          # Output: Mean prediction entropy: 1.2847
34          print(f"Entropy std: {all_entropies.std():.4f}")
35          # Output: Entropy std: 0.3214
36          print(f"Max entropy (most uncertain): {all_entropies.max():.4f}")
37          # Output: Max entropy (most uncertain): 2.1934
38          print(f"Min entropy (most certain): {all_entropies.min():.4f}")
39          # Output: Min entropy (most certain): 0.4756
40
41          return all_entropies
42
43  # Mutual information estimation (simplified)
44  def mutual_information_neural_estimation(x, y, hidden_dim=128):
45      """Neural estimation of mutual information"""
46      class MINENet(nn.Module):
47          def __init__(self, input_dim):
48              super().__init__()
49              self.net = nn.Sequential(
50                  nn.Linear(input_dim, hidden_dim),
51                  nn.ReLU(),
52                  nn.Linear(hidden_dim, hidden_dim),
53                  nn.ReLU(),
54                  nn.Linear(hidden_dim, 1)
55              )
56
57          def forward(self, x, y):
58              xy = torch.cat([x, y], dim=1)
59              return self.net(xy)
60
61      # This is a simplified version - full MINE requires more careful
          ↪   implementation
62      mine_net = MINENet(x.size(1) + y.size(1))
63      return mine_net
```

## 1.4 Optimization Theory

### 1.4.1 Gradient-Based Optimization

**Purpose:** Understanding optimization principles underlying deep learning training.
**Complex Example - Custom Optimizer:**

```python
1  class AdaptiveMomentumOptimizer(torch.optim.Optimizer):
2      """Custom optimizer implementing adaptive momentum"""
3
4      def __init__(self, params, lr=1e-3, beta1=0.9, beta2=0.999, eps=1e-8,
       ↪   weight_decay=0):
5          defaults = dict(lr=lr, beta1=beta1, beta2=beta2, eps=eps,
              ↪   weight_decay=weight_decay)
```

```python
            super().__init__(params, defaults)

    def step(self, closure=None):
        loss = None
        if closure is not None:
            loss = closure()

        for group in self.param_groups:
            for p in group['params']:
                if p.grad is None:
                    continue

                grad = p.grad.data
                if grad.is_sparse:
                    raise RuntimeError(
                        'Optimizer does not support sparse gradients')

                state = self.state[p]

                # State initialization
                if len(state) == 0:
                    state['step'] = 0
                    state['exp_avg'] = torch.zeros_like(p.data)
                    state['exp_avg_sq'] = torch.zeros_like(p.data)

                exp_avg, exp_avg_sq = state['exp_avg'], state['exp_avg_sq']
                beta1, beta2 = group['beta1'], group['beta2']

                state['step'] += 1

                # Weight decay
                if group['weight_decay'] != 0:
                    grad = grad.add(p.data, alpha=group['weight_decay'])

                # Exponential moving average of gradient values
                exp_avg.mul_(beta1).add_(grad, alpha=1 - beta1)

                # Exponential moving average of squared gradient values
                exp_avg_sq.mul_(beta2).addcmul_(grad, grad, value=1 - beta2)

                # Bias correction
                bias_correction1 = 1 - beta1 ** state['step']
                bias_correction2 = 1 - beta2 ** state['step']

                # Adaptive learning rate
                denom = (exp_avg_sq.sqrt() /
                    math.sqrt(bias_correction2)).add_(group['eps'])
                step_size = group['lr'] / bias_correction1
```

```python
52
53                     # Update parameters
54                     p.data.addcdiv_(exp_avg, denom, value=-step_size)
55
56             return loss
57
58 # Usage example with learning rate scheduling
59 def train_with_custom_optimizer(model, dataloader, epochs=10):
60     optimizer = AdaptiveMomentumOptimizer(model.parameters(), lr=0.001)
61     scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer,
    ↪   T_max=epochs)
62
63     for epoch in range(epochs):
64         for batch_x, batch_y in dataloader:
65             optimizer.zero_grad()
66
67             predictions = model(batch_x)
68             loss = F.cross_entropy(predictions, batch_y)
69
70             loss.backward()
71
72             # Gradient clipping
73             torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
74
75             optimizer.step()
76
77         scheduler.step()
78         print(f"Epoch {epoch}: LR = {scheduler.get_last_lr()[0]:.6f}")
79         # Output: Epoch 0: LR = 0.001000
```

# Chapter 2

# Preface

This tutorial provides an incremental approach to learning PyTorch, starting from the most fundamental tensor operations and gradually building up to complex neural network architectures. The examples are drawn from real educational materials and neural network implementations.

The progression follows a carefully designed path:

1. Fundamental tensor operations and data types

2. Mathematical operations and broadcasting

3. Automatic differentiation and gradients

4. Neural network building blocks

5. Optimization and training loops

6. Complete neural network architectures

Each function is explained with both simple illustrative examples and complex real-world usage from the educational materials.

# Chapter 3

# Fundamental Tensor Operations

## 3.1 Creating Tensors

### 3.1.1 torch.tensor()

**Purpose:** Creates a tensor from data (lists, arrays, scalars).
   **Syntax:** torch.tensor(data, dtype=None, device=None, requires_grad=False)
   **Simple Example:**

```python
import torch

# Creating tensors from different data types
scalar = torch.tensor(3.14)
vector = torch.tensor([1, 2, 3, 4])
matrix = torch.tensor([[1, 2], [3, 4]])

print(f"Scalar: {scalar}")
# Output: Scalar: tensor(3.1400)
print(f"Vector: {vector}")
# Output: Vector: tensor([1, 2, 3, 4])
print(f"Matrix: {matrix}")
# Output: Matrix: tensor([[1, 2],
#                         [3, 4]])

# PyTorch 2.x: Better device specification
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
tensor_on_device = torch.tensor([1, 2, 3], device=device)
print(f"Tensor on {device}: {tensor_on_device}")
# Output: Tensor on cpu: tensor([1, 2, 3])
```

**Complex Example from Educational Materials:**

```python
# From makemore bigram implementation
xs, ys = [], []
for w in words:
    chs = ['.'] + list(w) + ['.']
    for ch1, ch2 in zip(chs, chs[1:]):
```

```
6            ix1 = stoi[ch1]
7            ix2 = stoi[ch2]
8            xs.append(ix1)
9            ys.append(ix2)
10
11   xs = torch.tensor(xs)  # Input character indices
12   ys = torch.tensor(ys)  # Target character indices
13   print(f"Input shape: {xs.shape}, Target shape: {ys.shape}")
14   # Output: Input shape: torch.Size([32, 28, 28]), Target shape: torch.Size([32])
```

### 3.1.2   torch.zeros()

**Purpose:** Creates a tensor filled with zeros.
   **Syntax:** `torch.zeros(size, dtype=None, device=None, requires_grad=False)`
   **Simple Example:**

```
1    # Creating zero tensors of different shapes
2    zeros_1d = torch.zeros(5)
3    zeros_2d = torch.zeros(3, 4)
4    zeros_3d = torch.zeros(2, 3, 4)
5
6    print(f"1D zeros: {zeros_1d}")
7    # Output: 1D zeros: tensor([0., 0., 0., 0., 0.])
8    print(f"2D zeros shape: {zeros_2d.shape}")
9    # Output: 2D zeros shape: torch.Size([3, 4])
10   print(f"3D zeros shape: {zeros_3d.shape}")
11   # Output: 3D zeros shape: torch.Size([2, 3, 4])
```

**Complex Example from Educational Materials:**

```
1    # From makemore - creating bigram count matrix
2    N = torch.zeros((27, 27), dtype=torch.int32)
3
4    # Fill the matrix with bigram counts
5    for w in words:
6        chs = ['.'] + list(w) + ['.']
7        for ch1, ch2 in zip(chs, chs[1:]):
8            ix1 = stoi[ch1]
9            ix2 = stoi[ch2]
10           N[ix1, ix2] += 1
11
12   print(f"Bigram count matrix shape: {N.shape}")
13   # Output: Bigram count matrix shape: torch.Size([27, 27])
14   print(f"Total bigrams: {N.sum()}")
15   # Output: Total bigrams: tensor(32033)
```

### 3.1.3   torch.randn()

**Purpose:** Creates a tensor with random numbers from a normal distribution.

Syntax: `torch.randn(size, generator=None, dtype=None, device=None, requires_grad=False)`

**Simple Example:**

```python
# Creating random tensors
random_vector = torch.randn(5)
random_matrix = torch.randn(3, 3)

# Using a generator for reproducibility
g = torch.Generator().manual_seed(42)
reproducible_random = torch.randn(2, 3, generator=g)

print(f"Random vector: {random_vector}")
# Output: Random vector: tensor([-0.3420,  1.2341, -0.8765,  0.4321, -1.5432])
print(f"Random matrix:\n{random_matrix}")
# Output: Random matrix:
# tensor([[ 0.1234, -0.5678,  0.9876],
#         [-1.2345,  0.6789, -0.3456],
#         [ 0.7654, -0.9012,  1.3579]])

# PyTorch 2.x: Using device and dtype specifications
device = "cuda" if torch.cuda.is_available() else "cpu"
random_gpu = torch.randn(3, 3, device=device, dtype=torch.float32)
print(f"Random tensor on {device}: {random_gpu}")
# Output: Random tensor on cpu: tensor([[ 0.4567, -0.1234,  0.7890],
#                                       [-0.2345,  0.8901, -0.5678],
#                                       [ 0.3456, -0.7890,  0.1234]])
```

**Complex Example from Educational Materials:**

```python
# From makemore neural network initialization
g = torch.Generator().manual_seed(2147483647)
W = torch.randn((27, 27), generator=g, requires_grad=True)

# This creates the weight matrix for a neural network
# where each of 27 neurons receives 27 inputs
print(f"Weight matrix shape: {W.shape}")
# Output: Weight matrix shape: torch.Size([27, 27])
print(f"Requires gradient: {W.requires_grad}")
# Output: Requires gradient: True

# From Transformer initialization in makemore
config = ModelConfig(vocab_size=vocab_size, block_size=block_size,
                     n_layer=4, n_head=4, n_embd=64, n_embd2=64)
# Networks use randn internally for parameter initialization
```

### 3.1.4 torch.arange()

**Purpose:** Creates a tensor with a sequence of numbers.
   **Syntax:** `torch.arange(start, end, step=1, dtype=None, device=None)`
   **Simple Example:**

```python
# Creating sequences
seq1 = torch.arange(5)          # [0, 1, 2, 3, 4]
seq2 = torch.arange(1, 6)       # [1, 2, 3, 4, 5]
seq3 = torch.arange(0, 10, 2)   # [0, 2, 4, 6, 8]

print(f"Simple sequence: {seq1}")
# Output: Simple sequence: tensor([0, 1, 2, 3, 4])
print(f"Start-end sequence: {seq2}")
# Output: Start-end sequence: tensor([1, 2, 3, 4, 5])
print(f"With step: {seq3}")
# Output: With step: tensor([0, 2, 4, 6, 8])
```

**Complex Example from Educational Materials:**

```python
# From Transformer position embeddings
def forward(self, idx, targets=None):
    device = idx.device
    b, t = idx.size()
    assert t <= self.block_size

    # Create position indices for embeddings
    pos = torch.arange(0, t, dtype=torch.long, device=device).unsqueeze(0)

    # Get token and position embeddings
    tok_emb = self.transformer.wte(idx)  # (b, t, n_embd)
    pos_emb = self.transformer.wpe(pos)  # (1, t, n_embd)

    return tok_emb + pos_emb
```

## 3.2 Tensor Properties and Manipulation

### 3.2.1 Tensor.shape and Tensor.size()

**Purpose:** Get the dimensions of a tensor.
   **Simple Example:**

```python
tensor_2d = torch.randn(3, 4)
tensor_3d = torch.randn(2, 3, 4)

# Both .shape and .size() work
print(f"2D tensor shape: {tensor_2d.shape}")
```

```python
6    # Output: 2D tensor shape: torch.Size([3, 4])
7    print(f"2D tensor size: {tensor_2d.size()}")
8    # Output: 2D tensor size: torch.Size([3, 4])
9    print(f"3D tensor shape: {tensor_3d.shape}")
10   # Output: 3D tensor shape: torch.Size([2, 3, 4])
11
12   # Access specific dimensions
13   print(f"First dimension: {tensor_2d.shape[0]}")
14   # Output: First dimension: 3
15   print(f"Second dimension: {tensor_2d.size(1)}")
16   # Output: Second dimension: 4
```

**Complex Example from Educational Materials:**

```python
1    # From Transformer forward pass
2    def forward(self, x):
3        B, T, C = x.size()  # batch, sequence, embedding dimensions
4
5        # Split into query, key, value
6        q, k, v = self.c_attn(x).split(self.n_embd, dim=2)
7
8        # Reshape for multi-head attention
9        k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
10       q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
11       v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
12
13       print(f"Reshaped k: {k.shape}")  # (B, nh, T, hs)
14       # Output: Reshaped k: torch.Size([2, 8, 1024, 64])
15       return q, k, v
```

### 3.2.2 Tensor.view()

**Purpose:** Reshapes a tensor without changing its data.

**Simple Example:**

```python
1    # Original tensor
2    x = torch.arange(12)
3    print(f"Original: {x}")
4    # Output: Original: tensor([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
5
6    # Reshape to different dimensions
7    x_2d = x.view(3, 4)
8    x_3d = x.view(2, 2, 3)
9    x_flat = x_3d.view(-1)  # -1 means infer this dimension
10
11   print(f"2D view: {x_2d}")
```

```
12   # Output: 2D view: tensor([[ 0,  1,  2,  3],
13   #                           [ 4,  5,  6,  7],
14   #                           [ 8,  9, 10, 11]])
15   print(f"3D view: {x_3d}")
16   # Output: 3D view: tensor([[[ 0,  1,  2],
17   #                           [ 3,  4,  5]],
18   #                          [[ 6,  7,  8],
19   #                           [ 9, 10, 11]]])
20   print(f"Flattened: {x_flat}")
21   # Output: Flattened: tensor([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
```

**Complex Example from Educational Materials:**

```
1    # From Transformer multi-head attention
2    def forward(self, x):
3        B, T, C = x.size()
4
5        # Reshape for multi-head attention
6        k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
7
8        # After attention computation, reshape back
9        y = att @ v  # (B, nh, T, T) x (B, nh, T, hs) -> (B, nh, T, hs)
10       y = y.transpose(1, 2).contiguous().view(B, T, C)
11
12       # From loss computation - flatten for cross entropy
13       loss = F.cross_entropy(logits.view(-1, logits.size(-1)),
14                              targets.view(-1), ignore_index=-1)
15       return y
```

**Visual Guide to Tensor Reshaping:**

The following diagram shows how tensor.view() transforms data layout while preserving elements:

# Original: torch.arange(12)

[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]    **Shape: (12,)**

2D View: 3×4 Matrix    view(3, 4)

**3 rows**
**4 columns**
**Same 12**
**elements**    **Shape: (3, 4)**

3D View: 2×2×3 Tensor    view(2, 2, 3)

**Slice 0**
**2×3**
Elements
0-5

**Slice 1**
**2×3**
Elements
6-11    **Shape: (2, 2, 3)**

**Key Insights:**

- Elements maintain their order during reshaping

- Total number of elements must remain the same

- Use -1 to automatically infer one dimension

- Memory layout changes, but data is preserved

### 3.2.3 Tensor.unsqueeze() and Tensor.squeeze()

**Purpose:** Add or remove dimensions of size 1.
**Simple Example:**

```python
# Start with a 1D tensor
x = torch.tensor([1, 2, 3, 4])
print(f"Original shape: {x.shape}")
# Output: Original shape: torch.Size([4])


# Add dimensions
x_col = x.unsqueeze(1)      # Make column vector
x_row = x.unsqueeze(0)      # Make row vector
x_batch = x.unsqueeze(0).unsqueeze(0)  # Add batch and feature dims


print(f"Column vector: {x_col.shape}")
# Output: Column vector: torch.Size([4, 1])
print(f"Row vector: {x_row.shape}")
# Output: Row vector: torch.Size([1, 4])
print(f"With batch dim: {x_batch.shape}")
# Output: With batch dim: torch.Size([1, 1, 4])


# Remove dimensions of size 1
x_back = x_batch.squeeze()
print(f"After squeeze: {x_back.shape}")
# Output: After squeeze: torch.Size([4])
```

**Complex Example from Educational Materials:**

```python
# From position embeddings in Transformer
pos = torch.arange(0, t, dtype=torch.long, device=device).unsqueeze(0)
# Shape: (1, t) - adds batch dimension for broadcasting


# From sampling in makemore
xenc = F.one_hot(torch.tensor([ix]), num_classes=27).float()
# Creates one-hot vector for single character, unsqueeze for batch dim


# From keeping dimensions in softmax
P = (N+1).float()
P /= P.sum(1, keepdims=True)  # keepdims preserves dimension for broadcasting
```

**Visual Guide to Squeeze/Unsqueeze Operations:**

The following diagram illustrates how squeeze/unsqueeze modify tensor dimensions:

# Original Tensor: [1, 2, 3, 4]

1D Tensor
.unsqueeze(0)
**Shape: (4,)**

.unsqueeze(1)          .unsqueeze(0).unsqueeze(0)

**Column Vector    Row Vector    Batch Tensor**

4 rows 1 column

1 row 4 columns

1 batch 1 row 4 columns

**Shape: (4, 1)    Shape: (1, 4)    Shape: (1, 1, 4)**

.squeeze()

**After Squeeze**

Back to 1D

**Shape: (4,)**

**Broadcasting Visualization:**

Understanding how unsqueeze enables broadcasting:

## Vector A: (3,)    Vector B: (3,)

[1, 2, 3]          [4, 5, 6]

.unsqueeze(0)          .unsqueeze(1)

**A: (1, 3) Row    B: (3, 1) Column**

1 row, 3 cols

3 rows 1 col

**A + B broadcasts**

## Result: (3, 3) Matrix

3×3 Matrix Each element A[i] + B[j]

**Shape: (3, 3)**

**Key Concepts:**

- **unsqueeze(dim)**: Adds dimension of size 1 at specified position

- **squeeze()**: Removes all dimensions of size 1

- **squeeze(dim)**: Removes dimension of size 1 at specific position

- Essential for broadcasting operations between tensors

- Commonly used in CNN/RNN for batch dimension handling

# Chapter 4

# Mathematical Operations

## 4.1 Basic Arithmetic

### 4.1.1 Element-wise Operations

**Purpose:** Perform mathematical operations element by element.

**Simple Example:**

```python
a = torch.tensor([1, 2, 3, 4])
b = torch.tensor([2, 3, 4, 5])

# Basic arithmetic
add_result = a + b          # [3, 5, 7, 9]
sub_result = a - b          # [-1, -1, -1, -1]
mul_result = a * b          # [2, 6, 12, 20]
div_result = b / a          # [2.0, 1.5, 1.33, 1.25]
pow_result = a ** 2         # [1, 4, 9, 16]

print(f"Addition: {add_result}")
# Output: Addition: tensor([3, 5, 7, 9])
print(f"Multiplication: {mul_result}")
# Output: Multiplication: tensor([ 2,  6, 12, 20])
print(f"Power: {pow_result}")
# Output: Power: tensor([ 1,  4,  9, 16])
```

**Complex Example from Educational Materials:**

```python
# From makemore neural network forward pass
def forward(self, idx, targets=None):
    # Token and position embeddings
    tok_emb = self.transformer.wte(idx)  # (b, t, n_embd)
    pos_emb = self.transformer.wpe(pos)  # (1, t, n_embd)
    x = tok_emb + pos_emb  # Broadcasting addition

    # In regularization
    loss = -probs[torch.arange(num), ys].log().mean() + 0.01*(W**2).mean()
    #
```

```python
11      #                                          L2 regularization
12
13      # In gradient update
14      W.data += -50 * W.grad  # Element-wise multiplication and addition
```

## 4.1.2   Matrix Multiplication (@)

**Purpose:** Perform matrix multiplication (dot product).

   **Simple Example:**

```python
1   # 2D matrix multiplication
2   A = torch.randn(3, 4)
3   B = torch.randn(4, 5)
4   C = A @ B  # or torch.matmul(A, B)
5
6   print(f"A shape: {A.shape}")
7   # Output: A shape: torch.Size([3, 4])
8   print(f"B shape: {B.shape}")
9   # Output: B shape: torch.Size([4, 5])
10  print(f"C shape: {C.shape}")  # (3, 5)
11  # Output: C shape: torch.Size([3, 5])
12
13  # Vector-matrix multiplication
14  vec = torch.randn(3)
15  result = vec @ A  # (3,) @ (3, 4) -> (4,)
16  print(f"Vector-matrix result shape: {result.shape}")
17  # Output: Vector-matrix result shape: torch.Size([4])
```

   **Complex Example from Educational Materials:**

```python
1   # From neural network forward pass
2   def forward(self, x):
3       # Linear transformation
4       xenc = F.one_hot(xs, num_classes=27).float()
5       logits = xenc @ W  # (5, 27) @ (27, 27) -> (5, 27)
6
7       # Multi-head attention computation
8       att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
9       #      ^^^^^^^^^^^^^^^^^^^^^^^^^^^
10      #      Attention scores computation
11      y = att @ v  # (B, nh, T, T) @ (B, nh, T, hs) -> (B, nh, T, hs)
12
13      # In RNN cell
14      xh = torch.cat([xt, hprev], dim=1)
15      ht = F.tanh(self.xh_to_h(xh))  # Linear layer uses @ internally
16
```

```python
# PyTorch 2.x: Optimized matrix multiplication with torch.compile
@torch.compile
def optimized_matmul(A, B):
    return A @ B
```

**Visual Guide to Matrix Multiplication:**

Understanding how matrix dimensions work in multiplication:

## Matrix A: (3, 4) Matrix B: (4, 5) Result C: (3, 5)

| | | |
|---|---|---|
| **3 rows**<br>**4**<br>**columns** @ | **4 rows**<br>**5 columns** = | **3 rows**<br>**5 columns** |

**Inner dims match: 4 = 4**

**Outer dims: (3, 5)**

**Broadcasting in Matrix Operations:**

How PyTorch handles different dimensional operations:

## Vector @ Matrix

Vector (3,) @ Matrix (3,4) = Result (4,)

## Batch Matrix Multiplication

| Batch A<br>2<br>matrices<br>each<br>3×4 | (2, 4 @ Batch B<br>2<br>matrices<br>each<br>4×5 | 5) = Batch C<br>2<br>matrices<br>each<br>3×5 |
|---|---|---|

**Key Rules for Matrix Multiplication:**

- **Dimension compatibility**: Inner dimensions must match (A: m×n, B: n×p → C: m×p)

- **Batch operations**: PyTorch handles batched matrix multiplication automatically

- **Broadcasting**: Vector-matrix operations broadcast appropriately

- **@ operator**: Preferred over torch.matmul() for readability

- **Memory efficiency**: Use torch.compile for optimized operations

## 4.2 Activation Functions

### 4.2.1 torch.relu() and tensor.relu()

**Purpose:** Apply Rectified Linear Unit activation function.

**Simple Example:**

```python
import torch.nn.functional as F

x = torch.tensor([-2, -1, 0, 1, 2], dtype=torch.float)

# Three ways to apply ReLU
relu1 = torch.relu(x)
relu2 = F.relu(x)
relu3 = x.relu()  # Method on tensor

print(f"Input: {x}")
# Output: Input: tensor([-2., -1.,  0.,  1.,  2.])
print(f"ReLU output: {relu1}")
# Output: ReLU output: tensor([0., 0., 0., 1., 2.])

# ReLU zeros out negative values
negative_input = torch.randn(5)
positive_output = F.relu(negative_input)
print(f"Negative input: {negative_input}")
# Output: Negative input: tensor([-0.7324,  1.2356, -0.4567,  0.8901, -1.3245])
print(f"After ReLU: {positive_output}")
# Output: After ReLU: tensor([0.0000, 1.2356, 0.0000, 0.8901, 0.0000])
```

**Complex Example from Educational Materials:**

```python
# From micrograd Value class
def relu(self):
    out = Value(0 if self.data < 0 else self.data, (self,), 'ReLU')

    def _backward():
        self.grad += (out.data > 0) * out.grad
    out._backward = _backward
    return out

# From makemore generation with top-k sampling
if top_k is not None:
    v, _ = torch.topk(logits, top_k)
    # Apply ReLU-like behavior: set small values to -inf
    logits[logits < v[:, [-1]]] = -float('Inf')
```

### 4.2.2   torch.tanh()

**Purpose:** Apply hyperbolic tangent activation function.
    **Simple Example:**

```python
x = torch.linspace(-3, 3, 7)
tanh_output = torch.tanh(x)
```

```
3
4    print(f"Input: {x}")
5    # Output: Input: tensor([-3.0000, -2.0000, -1.0000,  0.0000,  1.0000,  2.0000,
     ↪  3.0000])
6    print(f"Tanh output: {tanh_output}")
7    # Output: Tanh output: tensor([-0.9951, -0.9640, -0.7616,  0.0000,  0.7616,
     ↪  0.9640,  0.9951])
8
9    # Tanh outputs are in range [-1, 1]
10   print(f"Min tanh: {tanh_output.min()}")
11   # Output: Min tanh: tensor(-0.9951)
12   print(f"Max tanh: {tanh_output.max()}")
13   # Output: Max tanh: tensor(0.9951)
```

**Complex Example from Educational Materials:**

```
1    # From RNN cell implementation
2    def forward(self, xt, hprev):
3        xh = torch.cat([xt, hprev], dim=1)
4        ht = F.tanh(self.xh_to_h(xh))  # Tanh activation for hidden state
5        return ht
6
7    # From GRU cell
8    def forward(self, xt, hprev):
9        # Calculate candidate hidden state
10       xhr = torch.cat([xt, hprev_reset], dim=1)
11       hbar = F.tanh(self.xh_to_hbar(xhr))
12
13       # Blend previous and candidate states
14       ht = (1 - z) * hprev + z * hbar
15       return ht
16
17   # From MLP layer
18   self.mlpf = lambda x: m.c_proj(F.tanh(m.c_fc(x)))
```

# Chapter 5

# Automatic Differentiation

## 5.1 requires_grad and Gradient Computation

### 5.1.1 requires_grad Parameter

**Purpose:** Enable automatic gradient computation for tensors.

**Computational Graph Visualization:**



$$z = x^2 + 2y + 1$$

**Simple Example:**

```python
# Create tensors that require gradients
x = torch.tensor(2.0, requires_grad=True)
y = torch.tensor(3.0, requires_grad=True)

# Perform operations
z = x**2 + 2*y + 1
print(f"z = {z}")
# Output: z = tensor(14., grad_fn=<AddBackward0>)

# Compute gradients
z.backward()
print(f"dz/dx = {x.grad}")   # Should be 2*x = 4
# Output: dz/dx = tensor(4.)
print(f"dz/dy = {y.grad}")   # Should be 2
# Output: dz/dy = tensor(2.)

# For tensors
```

```
18   W = torch.randn(2, 3, requires_grad=True)
19   print(f"W requires grad: {W.requires_grad}")
20   # Output: W requires grad: True
```

**Complex Example from Educational Materials:**

```
1    # From makemore neural network training
2    g = torch.Generator().manual_seed(2147483647)
3    W = torch.randn((27, 27), generator=g, requires_grad=True)
4
5    # Forward pass
6    xenc = F.one_hot(xs, num_classes=27).float()
7    logits = xenc @ W
8    counts = logits.exp()
9    probs = counts / counts.sum(1, keepdims=True)
10   loss = -probs[torch.arange(num), ys].log().mean() + 0.01*(W**2).mean()
11
12   # Backward pass
13   W.grad = None   # Clear previous gradients
14   loss.backward()   # Compute gradients
15
16   # Update parameters
17   W.data += -50 * W.grad
18   print(f"Gradient norm: {W.grad.norm()}")
19   # Output: Gradient norm: tensor(0.2847)
20
21   # PyTorch 2.x: Using autocast for mixed precision
22   with torch.autocast(device_type='cuda', enabled=torch.cuda.is_available()):
23       logits = xenc @ W
24       loss = F.cross_entropy(logits, ys)
```

### 5.1.2   tensor.backward()

**Purpose:** Compute gradients using backpropagation.
   **Simple Example:**

```
1    # Simple function: f(x, y) = x^2 + 3xy + y^2
2    x = torch.tensor(1.0, requires_grad=True)
3    y = torch.tensor(2.0, requires_grad=True)
4
5    # Forward pass
6    f = x**2 + 3*x*y + y**2
7    print(f"Function value: {f}")
8    # Output: Function value: tensor(11., grad_fn=<AddBackward0>)
9
10   # Backward pass
```

```
11   f.backward()
12
13   print(f"df/dx: {x.grad}")  # 2x + 3y = 2(1) + 3(2) = 8
14   # Output: df/dx: tensor(8.)
15   print(f"df/dy: {y.grad}")  # 3x + 2y = 3(1) + 2(2) = 7
16   # Output: df/dy: tensor(7.)
```

**Complex Example from Educational Materials:**

```
1    # From makemore training loop
2    for step in range(max_steps):
3        # Get batch
4        batch = batch_loader.next()
5        X, Y = [t.to(device) for t in batch]
6
7        # Forward pass
8        logits, loss = model(X, Y)
9
10       # Backward pass
11       model.zero_grad(set_to_none=True)  # Clear gradients
12       loss.backward()                    # Compute gradients
13       optimizer.step()                   # Update parameters
14
15       if step % 10 == 0:
16           print(f"step {step} | loss {loss.item():.4f}")
17           # Output: step 0 | loss 2.4567
18
19   # PyTorch 2.x: Using torch.compile for optimization
20   model = torch.compile(model)  # Faster execution
21
22   # PyTorch 2.x: Better mixed precision training
23   scaler = torch.cuda.amp.GradScaler()
24   with torch.autocast(device_type='cuda'):
25       logits, loss = model(X, Y)
26
27   scaler.scale(loss).backward()
28   scaler.step(optimizer)
29   scaler.update()
30
31   # From micrograd implementation
32   def backward(self):
33       # Build topological order
34       topo = []
35       visited = set()
36       def build_topo(v):
37           if v not in visited:
38               visited.add(v)
```

```python
39                for child in v._prev:
40                    build_topo(child)
41                topo.append(v)
42        build_topo(self)
43
44        # Apply chain rule
45        self.grad = 1
46        for v in reversed(topo):
47            v._backward()
```

# Chapter 6

# Convolutional Neural Networks

## 6.1 Convolutional Layers

### 6.1.1 torch.nn.Conv2d

**Purpose:** Applies 2D convolution for image processing and feature extraction.

Syntax: `nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0)`

**Convolution Operation Visualization:**

**Input Feature Map**

**3×3 Kernel** **Output Feature Map**

Padding

5×5×3

*

3×3×3

3×3×64

Stride = 1

**Simple Example:**

```python
import torch.nn as nn

# Basic convolution layer
conv = nn.Conv2d(
    in_channels=3,      # RGB input
    out_channels=64,    # 64 feature maps
    kernel_size=3,      # 3x3 kernel
    stride=1,           # Stride of 1
    padding=1           # Padding to keep size
)

# Input: batch_size=32, channels=3, height=224, width=224
x = torch.randn(32, 3, 224, 224)
output = conv(x)
print(f"Input shape: {x.shape}")      # torch.Size([32, 3, 224, 224])
print(f"Output shape: {output.shape}") # torch.Size([32, 64, 224, 224])
```

```python
18  # Access layer parameters
19  print(f"Weight shape: {conv.weight.shape}")  # torch.Size([64, 3, 3, 3])
20  print(f"Bias shape: {conv.bias.shape}")      # torch.Size([64])
```

**Complex Example - CNN Architecture:**

```python
1   class CNN_Classifier(nn.Module):
2       def __init__(self, num_classes=10):
3           super().__init__()
4
5           # Feature extraction layers
6           self.features = nn.Sequential(
7               # First conv block
8               nn.Conv2d(3, 64, kernel_size=3, padding=1),
9               nn.BatchNorm2d(64),
10              nn.ReLU(inplace=True),
11              nn.MaxPool2d(kernel_size=2, stride=2),
12
13              # Second conv block
14              nn.Conv2d(64, 128, kernel_size=3, padding=1),
15              nn.BatchNorm2d(128),
16              nn.ReLU(inplace=True),
17              nn.MaxPool2d(kernel_size=2, stride=2),
18
19              # Third conv block
20              nn.Conv2d(128, 256, kernel_size=3, padding=1),
21              nn.BatchNorm2d(256),
22              nn.ReLU(inplace=True),
23              nn.MaxPool2d(kernel_size=2, stride=2),
24          )
25
26          # Classifier
27          self.classifier = nn.Sequential(
28              nn.AdaptiveAvgPool2d((7, 7)),
29              nn.Flatten(),
30              nn.Linear(256 * 7 * 7, 512),
31              nn.ReLU(inplace=True),
32              nn.Dropout(0.5),
33              nn.Linear(512, num_classes)
34          )
35
36      def forward(self, x):
37          x = self.features(x)
38          x = self.classifier(x)
39          return x
40
41  # Usage
```

```
42   model = CNN_Classifier(num_classes=1000)
43   input_tensor = torch.randn(16, 3, 224, 224)
44   output = model(input_tensor)
45   print(f"Output shape: {output.shape}")  # torch.Size([16, 1000])
```

**CNN Architecture Visualization:**



## 6.2 Pooling Layers

### 6.2.1 torch.nn.MaxPool2d

**Purpose:** Applies max pooling for downsampling and translation invariance.

**Max Pooling Operation Visualization:**



Kernel: 2×2, Stride: 2

**Simple Example:**

```python
1    # Max pooling layer
2    maxpool = nn.MaxPool2d(
3        kernel_size=2,    # 2x2 pooling window
4        stride=2,         # Non-overlapping windows
5        padding=0         # No padding
6    )
7
8    # Input: 64 feature maps of size 56x56
9    x = torch.randn(32, 64, 56, 56)
10   output = maxpool(x)
11   print(f"Input shape: {x.shape}")      # torch.Size([32, 64, 56, 56])
12   print(f"Output shape: {output.shape}") # torch.Size([32, 64, 28, 28])
13
14   # Average pooling
```

```python
15    avgpool = nn.AvgPool2d(kernel_size=2, stride=2)
16    avg_output = avgpool(x)
17    print(f"AvgPool output: {avg_output.shape}")  # torch.Size([32, 64, 28, 28])
```

**Complex Example - Adaptive Pooling:**

```python
1    # Adaptive pooling - always produces fixed output size
2    adaptive_avg = nn.AdaptiveAvgPool2d((7, 7))  # Always 7x7 output
3    adaptive_max = nn.AdaptiveMaxPool2d((1, 1))  # Global pooling
4
5    # Different input sizes
6    inputs = [
7        torch.randn(1, 256, 14, 14),  # Small feature map
8        torch.randn(1, 256, 28, 28),  # Medium feature map
9        torch.randn(1, 256, 56, 56),  # Large feature map
10   ]
11
12   for i, input_tensor in enumerate(inputs):
13       # Adaptive average pooling
14       avg_out = adaptive_avg(input_tensor)
15       max_out = adaptive_max(input_tensor)
16
17       print(f"Input {i+1}: {input_tensor.shape}")
18       print(f"  Adaptive Avg: {avg_out.shape}")      # Always (1, 256, 7, 7)
19       print(f"  Adaptive Max: {max_out.shape}")      # Always (1, 256, 1, 1)
20
21   # Global Average Pooling (common in modern architectures)
22   class GlobalAvgPool(nn.Module):
23       def forward(self, x):
24           # x: (batch_size, channels, height, width)
25           return F.adaptive_avg_pool2d(x, (1, 1)).view(x.size(0), -1)
26
27   gap = GlobalAvgPool()
28   x = torch.randn(32, 512, 7, 7)
29   output = gap(x)
30   print(f"Global pooling output: {output.shape}")  # torch.Size([32, 512])
```

## 6.3  Activation Functions

### 6.3.1  torch.nn.GELU

**Purpose:** Gaussian Error Linear Unit - modern activation function used in transformers.
   **Simple Example:**

```python
1    # GELU activation (used in BERT, GPT)
2    gelu = nn.GELU()
```

```
3    x = torch.randn(5)
4    output = gelu(x)
5
6    print(f"Input: {x}")
7    # Output: Input: tensor([-2.0000, -1.0000,  0.0000,  1.0000,  2.0000])
8    print(f"GELU output: {output}")
9    # Output: GELU output: tensor([-0.0455, -0.1588,  0.0000,  0.8413,  1.9545])
10
11   # Compare with ReLU
12   relu = nn.ReLU()
13   relu_output = relu(x)
14   print(f"ReLU output: {relu_output}")
15   # Output: ReLU output: tensor([0.0000, 0.0000, 0.0000, 1.0000, 2.0000])
16
17   # GELU is smoother than ReLU, allowing small negative values
```

**Complex Example - Modern Activation Functions:**

```
1    # Comparison of modern activation functions
2    activations = {
3        'ReLU': nn.ReLU(),
4        'GELU': nn.GELU(),
5        'SiLU (Swish)': nn.SiLU(),
6        'Mish': nn.Mish(),
7        'LeakyReLU': nn.LeakyReLU(0.1)
8    }
9
10   x = torch.linspace(-3, 3, 100)
11
12   # Test all activations
13   for name, activation in activations.items():
14       y = activation(x)
15       print(f"{name}: min={y.min():.3f}, max={y.max():.3f}")
16       # Output: ReLU: min=0.000, max=2.000
17       # Output: GELU: min=-0.046, max=1.955
18       # Output: Tanh: min=-0.995, max=0.995
19
20   # Modern MLP with GELU
21   class ModernMLP(nn.Module):
22       def __init__(self, input_dim, hidden_dim, output_dim):
23           super().__init__()
24           self.layers = nn.Sequential(
25               nn.Linear(input_dim, hidden_dim),
26               nn.GELU(),                        # Modern activation
27               nn.LayerNorm(hidden_dim),         # Layer normalization
28               nn.Dropout(0.1),
29               nn.Linear(hidden_dim, hidden_dim),
```

```python
30            nn.GELU(),
31            nn.LayerNorm(hidden_dim),
32            nn.Dropout(0.1),
33            nn.Linear(hidden_dim, output_dim)
34        )
35
36    def forward(self, x):
37        return self.layers(x)
```

# Chapter 7

# Neural Network Building Blocks

## 7.1 Linear Layers

### 7.1.1 torch.nn.Linear

**Purpose:** Applies a linear transformation: $y = xW^T + b$.

**Neural Network Architecture Visualization:**

**Input Layer**

**Hidden Layer 1**

**Hidden Layer 2**

**Output Layer**

$x_5$

$x_4$    $W_1, b_1$

$x_3$        $W_2, b_2$

$x_2$           $W_3, b_3$   $y_{3.5}$

$x_1$             $y_{2.5}$

(5 features)   (4 units)   (3 units)   (2 classes)

ReLU   ReLU   Softmax

**Simple Example:**

```python
import torch.nn as nn

# Create a linear layer
linear = nn.Linear(in_features=5, out_features=3)
print(f"Weight shape: {linear.weight.shape}")  # (3, 5)
# Output: Weight shape: torch.Size([3, 5])
print(f"Bias shape: {linear.bias.shape}")       # (3,)
# Output: Bias shape: torch.Size([3])

# Forward pass
x = torch.randn(2, 5)  # Batch of 2 samples, 5 features each
y = linear(x)
print(f"Input shape: {x.shape}")   # (2, 5)
# Output: Input shape: torch.Size([2, 5])
print(f"Output shape: {y.shape}")  # (2, 3)
# Output: Output shape: torch.Size([2, 3])
```

```
17
18    # PyTorch 2.x: Using device and dtype initialization
19    device = "cuda" if torch.cuda.is_available() else "cpu"
20    linear_gpu = nn.Linear(5, 3, device=device, dtype=torch.float16)
21    print(f"Layer on {device} with dtype {linear_gpu.weight.dtype}")
22    # Output: Layer on cpu with dtype torch.float32
```

**Complex Example from Educational Materials:**

```
1     # From Transformer implementation
2     class CausalSelfAttention(nn.Module):
3         def __init__(self, config):
4             super().__init__()
5             # Key, query, value projections for all heads
6             self.c_attn = nn.Linear(config.n_embd, 3 * config.n_embd)
7             # Output projection
8             self.c_proj = nn.Linear(config.n_embd, config.n_embd)
9
10        def forward(self, x):
11            # Apply linear transformations
12            q, k, v = self.c_attn(x).split(self.n_embd, dim=2)
13            y = self.c_proj(y)
14            return y
15
16    # From MLP implementation
17    class MLP(nn.Module):
18        def __init__(self, config):
19            super().__init__()
20            self.mlp = nn.Sequential(
21                nn.Linear(self.block_size * config.n_embd, config.n_embd2),
22                nn.Tanh(),
23                nn.Linear(config.n_embd2, self.vocab_size)
24            )
```

## 7.2  Embedding Layers

### 7.2.1  torch.nn.Embedding

**Purpose:** Creates learnable lookup tables for discrete tokens.

**Simple Example:**

```
1     # Create embedding layer
2     vocab_size = 1000
3     embedding_dim = 128
4     embedding = nn.Embedding(vocab_size, embedding_dim)
5
```

```
6    # Input: token indices
7    tokens = torch.tensor([1, 5, 23, 100])
8    embedded = embedding(tokens)
9
10   print(f"Token indices: {tokens}")
11   # Output: Token indices: tensor([  1,    5,   23, 100])
12   print(f"Embedded shape: {embedded.shape}")  # (4, 128)
13   # Output: Embedded shape: torch.Size([4, 128])
14   print(f"Each token -> {embedding_dim}D vector")
15   # Output: Each token -> 128D vector
16
17   # Batch processing
18   batch_tokens = torch.tensor([[1, 5, 23], [45, 67, 89]])
19   batch_embedded = embedding(batch_tokens)
20   print(f"Batch embedded shape: {batch_embedded.shape}")  # (2, 3, 128)
21   # Output: Batch embedded shape: torch.Size([2, 3, 128])
```

**Complex Example from Educational Materials:**

```
1    # From Transformer language model
2    class Transformer(nn.Module):
3        def __init__(self, config):
4            super().__init__()
5            self.transformer = nn.ModuleDict(dict(
6                wte = nn.Embedding(config.vocab_size, config.n_embd),    # Token
                  ↪ embeddings
7                wpe = nn.Embedding(config.block_size, config.n_embd),    # Position
                  ↪ embeddings
8                h = nn.ModuleList([Block(config) for _ in range(config.n_layer)]),
9                ln_f = nn.LayerNorm(config.n_embd),
10           ))
11
12       def forward(self, idx, targets=None):
13           b, t = idx.size()
14           pos = torch.arange(0, t, dtype=torch.long, device=device).unsqueeze(0)
15
16           # Get embeddings
17           tok_emb = self.transformer.wte(idx)  # Token embeddings
18           pos_emb = self.transformer.wpe(pos)  # Position embeddings
19           x = tok_emb + pos_emb  # Combine both
20
21           return x
22
23   # From MLP model
24   self.wte = nn.Embedding(config.vocab_size + 1, config.n_embd)
25   # +1 for special <BLANK> token
```

## 7.3   Normalization Layers

### 7.3.1   torch.nn.Dropout

**Purpose:** Applies dropout regularization to prevent overfitting.
**Simple Example:**

```python
# Dropout layer
dropout = nn.Dropout(p=0.5)  # Drop 50% of neurons during training

# Input tensor
x = torch.randn(32, 128)
print(f"Input: {x[0, :10]}")  # First 10 values of first sample
# Output: Input: tensor([ 0.7342, -1.2456,  0.9876, -0.3421,  1.5643, -0.8765,
↪  0.2134, -1.6789,  0.4321, -0.7654])

# During training (dropout active)
model.train()
dropped = dropout(x)
print(f"Dropped: {dropped[0, :10]}")  # Some values are zero
# Output: Dropped: tensor([ 1.4684,  0.0000,  1.9752, -0.6842,  0.0000, -1.7530,
↪  0.4268,  0.0000,  0.8642,  0.0000])

# During evaluation (dropout inactive)
model.eval()
eval_output = dropout(x)
print(f"Eval: {eval_output[0, :10]}")  # All values preserved
# Output: Eval: tensor([ 0.7342, -1.2456,  0.9876, -0.3421,  1.5643, -0.8765,
↪  0.2134, -1.6789,  0.4321, -0.7654])
```

**Complex Example - Different Dropout Types:**

```python
class DropoutComparison(nn.Module):
    def __init__(self):
        super().__init__()
        # Standard dropout
        self.dropout = nn.Dropout(0.3)

        # Dropout for convolutional layers
        self.dropout2d = nn.Dropout2d(0.25)  # Drops entire feature maps

        # Alpha dropout (for SELU activation)
        self.alpha_dropout = nn.AlphaDropout(0.3)

        # Feature alpha dropout
        self.feature_alpha_dropout = nn.FeatureAlphaDropout(0.3)

    def forward(self, x_linear, x_conv):
```

```
17            # Linear layer dropout
18            x_linear = self.dropout(x_linear)
19
20            # Convolutional dropout (drops entire channels)
21            x_conv = self.dropout2d(x_conv)
22
23            return x_linear, x_conv
24
25    # Example usage in CNN
26    class RegularizedCNN(nn.Module):
27        def __init__(self, num_classes):
28            super().__init__()
29            self.features = nn.Sequential(
30                nn.Conv2d(3, 64, 3, padding=1),
31                nn.ReLU(),
32                nn.Dropout2d(0.1),   # Spatial dropout
33
34                nn.Conv2d(64, 128, 3, padding=1),
35                nn.ReLU(),
36                nn.MaxPool2d(2),
37                nn.Dropout2d(0.2),
38            )
39
40            self.classifier = nn.Sequential(
41                nn.Linear(128 * 14 * 14, 512),
42                nn.ReLU(),
43                nn.Dropout(0.5),   # Standard dropout
44                nn.Linear(512, num_classes)
45            )
46
47        def forward(self, x):
48            x = self.features(x)
49            x = x.view(x.size(0), -1)
50            x = self.classifier(x)
51            return x
```

### 7.3.2   torch.nn.BatchNorm2d

**Purpose:** Applies batch normalization for faster training and regularization.

**Simple Example:**

```
1    # Batch normalization for 2D inputs (after Conv2d)
2    batch_norm = nn.BatchNorm2d(64)   # 64 feature channels
3
4    # Input: (batch_size, channels, height, width)
5    x = torch.randn(32, 64, 28, 28)
6    normalized = batch_norm(x)
```

```python
7
8    print(f"Input shape: {x.shape}")
9    # Output: Input shape: torch.Size([32, 64, 28, 28])
10   print(f"Output shape: {normalized.shape}")
11   # Output: Output shape: torch.Size([32, 64, 28, 28])
12
13   # Check statistics
14   print(f"Mean per channel: {normalized.mean(dim=[0,2,3])}")   # Should be ~0
15   # Output: Mean per channel: tensor([-0.0012,  0.0023, -0.0045, ...,  0.0018])
16   print(f"Std per channel: {normalized.std(dim=[0,2,3])}")      # Should be ~1
17   # Output: Std per channel: tensor([0.9987, 1.0034, 0.9956, ..., 1.0012])
18
19   # Access learned parameters
20   print(f"Gamma (scale): {batch_norm.weight.shape}")   # (64,)
21   # Output: Gamma (scale): torch.Size([64])
22   print(f"Beta (shift): {batch_norm.bias.shape}")       # (64,)
23   # Output: Beta (shift): torch.Size([64])
```

**Complex Example - Normalization Comparison:**

```python
1    class NormalizationComparison(nn.Module):
2        def __init__(self, channels, height, width):
3            super().__init__()
4            # Different normalization techniques
5            self.batch_norm = nn.BatchNorm2d(channels)
6            self.layer_norm = nn.LayerNorm([channels, height, width])
7            self.instance_norm = nn.InstanceNorm2d(channels)
8            self.group_norm = nn.GroupNorm(8, channels)  # 8 groups
9
10       def forward(self, x):
11           # x shape: (batch, channels, height, width)
12
13           # Batch normalization: normalize across batch dimension
14           bn_out = self.batch_norm(x)
15
16           # Layer normalization: normalize across channel, height, width
17           ln_out = self.layer_norm(x)
18
19           # Instance normalization: normalize per instance per channel
20           in_out = self.instance_norm(x)
21
22           # Group normalization: normalize within groups of channels
23           gn_out = self.group_norm(x)
24
25           return {
26               'batch_norm': bn_out,
27               'layer_norm': ln_out,
```

```
28              'instance_norm': in_out,
29              'group_norm': gn_out
30          }
31
32      # Modern CNN block with proper normalization
33      class ModernConvBlock(nn.Module):
34          def __init__(self, in_channels, out_channels, use_residual=False):
35              super().__init__()
36              self.use_residual = use_residual
37
38              self.conv1 = nn.Conv2d(in_channels, out_channels, 3, padding=1,
                ↪  bias=False)
39              self.bn1 = nn.BatchNorm2d(out_channels)
40              self.conv2 = nn.Conv2d(out_channels, out_channels, 3, padding=1,
                ↪  bias=False)
41              self.bn2 = nn.BatchNorm2d(out_channels)
42
43              # Residual connection
44              if use_residual and in_channels != out_channels:
45                  self.residual = nn.Conv2d(in_channels, out_channels, 1, bias=False)
46              else:
47                  self.residual = nn.Identity()
48
49          def forward(self, x):
50              identity = x
51
52              out = F.relu(self.bn1(self.conv1(x)))
53              out = self.bn2(self.conv2(out))
54
55              if self.use_residual:
56                  out += self.residual(identity)
57
58              return F.relu(out)
```

### 7.3.3  torch.nn.LayerNorm

**Purpose:** Applies layer normalization to stabilize training.

**Simple Example:**

```
1   # Create layer norm
2   layer_norm = nn.LayerNorm(4)  # Normalize over last dimension
3
4   # Input tensor
5   x = torch.randn(2, 3, 4)  # (batch, sequence, features)
6   normalized = layer_norm(x)
7
8   print(f"Input shape: {x.shape}")
```

```python
9    # Output: Input shape: torch.Size([2, 3, 4])
10   print(f"Output shape: {normalized.shape}")
11   # Output: Output shape: torch.Size([2, 3, 4])
12
13   # Check normalization: mean approximately 0, std approximately 1 for last
     ↪  dimension
14   print(f"Mean along last dim: {normalized.mean(dim=-1)}")
15   # Output: Mean along last dim: tensor([[-0.0000, -0.0000,  0.0000],
16   #                                      [ 0.0000,  0.0000, -0.0000]])
17   print(f"Std along last dim: {normalized.std(dim=-1)}")
18   # Output: Std along last dim: tensor([[1.0000, 1.0000, 1.0000],
19   #                                     [1.0000, 1.0000, 1.0000]])
```

**Complex Example from Educational Materials:**

```python
1    # From Transformer block
2    class Block(nn.Module):
3        def __init__(self, config):
4            super().__init__()
5            self.ln_1 = nn.LayerNorm(config.n_embd)  # Pre-attention norm
6            self.attn = CausalSelfAttention(config)
7            self.ln_2 = nn.LayerNorm(config.n_embd)  # Pre-MLP norm
8            self.mlp = nn.ModuleDict(dict(
9                c_fc    = nn.Linear(config.n_embd, 4 * config.n_embd),
10               c_proj  = nn.Linear(4 * config.n_embd, config.n_embd),
11               act     = NewGELU(),
12           ))
13
14       def forward(self, x):
15           # Pre-norm architecture
16           x = x + self.attn(self.ln_1(x))     # Residual + attention
17           x = x + self.mlpf(self.ln_2(x))     # Residual + MLP
18           return x
19
20   # Final layer norm before output
21   self.transformer = nn.ModuleDict(dict(
22       wte = nn.Embedding(config.vocab_size, config.n_embd),
23       wpe = nn.Embedding(config.block_size, config.n_embd),
24       h = nn.ModuleList([Block(config) for _ in range(config.n_layer)]),
25       ln_f = nn.LayerNorm(config.n_embd),  # Final layer norm
26   ))
```

# Chapter 8

# Essential Deep Learning Utilities

## 8.1 Gradient Clipping

### 8.1.1 torch.nn.utils.clip_grad_norm_()

**Purpose:** Clips gradient norm of parameters to prevent gradient explosion.

Syntax: `torch.nn.utils.clip_grad_norm_(parameters, max_norm, norm_type=2.0)`

**Simple Example:**

```python
import torch
import torch.nn as nn
import torch.nn.utils as utils

# Simple model
model = nn.Sequential(
    nn.Linear(10, 50),
    nn.ReLU(),
    nn.Linear(50, 1)
)

# Training step with gradient clipping
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
criterion = nn.MSELoss()

x = torch.randn(32, 10)
y = torch.randn(32, 1)

# Forward pass
output = model(x)
loss = criterion(output, y)

# Backward pass with gradient clipping
optimizer.zero_grad()
loss.backward()

# Clip gradients before optimizer step
grad_norm = utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
```

```
29    print(f"Gradient norm before clipping: {grad_norm}")
30    # Output: Gradient norm before clipping: tensor(2.4567)
31
32    optimizer.step()
```

**Complex Example - RNN Training:**

```
1    # Training loop with gradient clipping for RNN
2    class LSTM_Model(nn.Module):
3        def __init__(self, vocab_size, embed_size, hidden_size, num_layers):
4            super().__init__()
5            self.embedding = nn.Embedding(vocab_size, embed_size)
6            self.lstm = nn.LSTM(embed_size, hidden_size, num_layers, batch_first=True)
7            self.fc = nn.Linear(hidden_size, vocab_size)
8
9        def forward(self, x, hidden=None):
10            embedded = self.embedding(x)
11            lstm_out, hidden = self.lstm(embedded, hidden)
12            output = self.fc(lstm_out)
13            return output, hidden
14
15    model = LSTM_Model(vocab_size=10000, embed_size=256, hidden_size=512,
    ↪    num_layers=2)
16    optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
17
18    for epoch in range(num_epochs):
19        for batch in dataloader:
20            input_ids, target_ids = batch
21
22            # Forward pass
23            output, _ = model(input_ids)
24            loss = F.cross_entropy(output.view(-1, vocab_size), target_ids.view(-1))
25
26            # Backward pass with gradient clipping
27            optimizer.zero_grad()
28            loss.backward()
29
30            # Essential for RNN training - prevents exploding gradients
31            torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=0.5)
32
33            optimizer.step()
```

### 8.1.2   torch.nn.utils.clip_grad_value_()

**Purpose:** Clips gradients at a specified value.

   **Simple Example:**

```
1   # Value-based gradient clipping
2   optimizer.zero_grad()
3   loss.backward()
4
5   # Clip individual gradient values to [-0.5, 0.5]
6   torch.nn.utils.clip_grad_value_(model.parameters(), clip_value=0.5)
7
8   optimizer.step()
```

## 8.2   Model Utilities

### 8.2.1   torch.save() and torch.load()

**Purpose:** Save and load models, optimizers, and training state.

**Simple Example:**

```
1   # Save model state dict (recommended)
2   torch.save(model.state_dict(), 'model_weights.pth')
3
4   # Load model state dict
5   model = MyModel()
6   model.load_state_dict(torch.load('model_weights.pth'))
7   model.eval()
8
9   # Save entire model (less flexible)
10  torch.save(model, 'complete_model.pth')
11  loaded_model = torch.load('complete_model.pth')
```

**Complex Example - Training Checkpoint:**

```
1   # Save complete training state
2   def save_checkpoint(model, optimizer, scheduler, epoch, loss, filepath):
3       checkpoint = {
4           'epoch': epoch,
5           'model_state_dict': model.state_dict(),
6           'optimizer_state_dict': optimizer.state_dict(),
7           'scheduler_state_dict': scheduler.state_dict(),
8           'loss': loss,
9           'model_config': model.config  # Save model configuration
10      }
11      torch.save(checkpoint, filepath)
12
13  # Load complete training state
14  def load_checkpoint(filepath, model, optimizer=None, scheduler=None):
15      checkpoint = torch.load(filepath, map_location='cpu')
16
```

```python
17        model.load_state_dict(checkpoint['model_state_dict'])
18
19        if optimizer:
20            optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
21
22        if scheduler:
23            scheduler.load_state_dict(checkpoint['scheduler_state_dict'])
24
25        return checkpoint['epoch'], checkpoint['loss']
26
27    # Usage in training loop
28    if epoch % save_interval == 0:
29        save_checkpoint(model, optimizer, scheduler, epoch, loss,
30                        f'checkpoint_epoch_{epoch}.pth')
31
32    # Resume training
33    epoch_start, prev_loss = load_checkpoint('checkpoint_epoch_100.pth',
34                                             model, optimizer, scheduler)
```

## 8.3   Parameter Initialization

### 8.3.1   torch.nn.init Functions

**Purpose:** Initialize model parameters with specific distributions.

   **Simple Example:**

```python
1   import torch.nn.init as init
2
3   # Initialize a linear layer
4   layer = nn.Linear(100, 50)
5
6   # Xavier/Glorot initialization
7   init.xavier_uniform_(layer.weight)
8   init.zeros_(layer.bias)
9
10  # Kaiming/He initialization (good for ReLU)
11  init.kaiming_normal_(layer.weight, mode='fan_out', nonlinearity='relu')
12
13  # Normal initialization
14  init.normal_(layer.weight, mean=0, std=0.01)
15
16  # Constant initialization
17  init.constant_(layer.bias, 0)
```

   **Complex Example - Custom Model Initialization:**

```python
class CustomCNN(nn.Module):
    def __init__(self, num_classes):
        super().__init__()
        self.conv1 = nn.Conv2d(3, 64, kernel_size=3, padding=1)
        self.conv2 = nn.Conv2d(64, 128, kernel_size=3, padding=1)
        self.fc1 = nn.Linear(128 * 8 * 8, 512)
        self.fc2 = nn.Linear(512, num_classes)

        # Custom initialization
        self._initialize_weights()

    def _initialize_weights(self):
        for module in self.modules():
            if isinstance(module, nn.Conv2d):
                # Kaiming initialization for conv layers
                nn.init.kaiming_normal_(module.weight, mode='fan_out',
                                        nonlinearity='relu')
                if module.bias is not None:
                    nn.init.constant_(module.bias, 0)
            elif isinstance(module, nn.Linear):
                # Xavier initialization for linear layers
                nn.init.xavier_normal_(module.weight)
                nn.init.constant_(module.bias, 0)

# Alternative: Initialize after model creation
def init_weights(module):
    if isinstance(module, nn.Linear):
        torch.nn.init.xavier_uniform_(module.weight)
        module.bias.data.fill_(0.01)
    elif isinstance(module, nn.Conv2d):
        torch.nn.init.kaiming_uniform_(module.weight)

model = CustomCNN(num_classes=10)
model.apply(init_weights)  # Apply to all modules
```

## 8.4 Data Utilities

### 8.4.1 torch.utils.data.DataLoader

**Purpose:** Efficient data loading with batching, shuffling, and multiprocessing.
Simple Example:

```python
from torch.utils.data import DataLoader, TensorDataset

# Create dataset
x_data = torch.randn(1000, 10)
```

```python
5    y_data = torch.randn(1000, 1)
6    dataset = TensorDataset(x_data, y_data)
7
8    # Create dataloader
9    dataloader = DataLoader(
10       dataset,
11       batch_size=32,
12       shuffle=True,
13       num_workers=4,
14       pin_memory=True  # Faster GPU transfer
15   )
16
17   # Training loop
18   for batch_idx, (data, targets) in enumerate(dataloader):
19       # Move to GPU
20       data = data.to(device)
21       targets = targets.to(device)
22
23       # Training step
24       outputs = model(data)
25       loss = criterion(outputs, targets)
```

**Complex Example - Custom Dataset:**

```python
1    from torch.utils.data import Dataset, DataLoader
2    from PIL import Image
3    import torchvision.transforms as transforms
4
5    class CustomImageDataset(Dataset):
6        def __init__(self, image_paths, labels, transform=None):
7            self.image_paths = image_paths
8            self.labels = labels
9            self.transform = transform
10
11       def __len__(self):
12           return len(self.image_paths)
13
14       def __getitem__(self, idx):
15           # Load image
16           image = Image.open(self.image_paths[idx])
17           label = self.labels[idx]
18
19           # Apply transforms
20           if self.transform:
21               image = self.transform(image)
22
23           return image, label
```

```
24
25    # Define transforms
26    transform_train = transforms.Compose([
27        transforms.RandomResizedCrop(224),
28        transforms.RandomHorizontalFlip(),
29        transforms.ColorJitter(brightness=0.4, contrast=0.4, saturation=0.4),
30        transforms.ToTensor(),
31        transforms.Normalize(mean=[0.485, 0.456, 0.406],
32                            std=[0.229, 0.224, 0.225])
33    ])
34
35    # Create dataset and dataloader
36    train_dataset = CustomImageDataset(train_paths, train_labels, transform_train)
37    train_loader = DataLoader(
38        train_dataset,
39        batch_size=64,
40        shuffle=True,
41        num_workers=8,
42        pin_memory=True,
43        persistent_workers=True,  # PyTorch 2.x feature
44        prefetch_factor=2
45    )
```

### 8.4.2   torch.stack() and torch.cat()

**Purpose:** Combine tensors along different dimensions.

**Simple Example:**

```
1     # torch.cat - concatenate along existing dimension
2     a = torch.tensor([[1, 2], [3, 4]])
3     b = torch.tensor([[5, 6], [7, 8]])
4
5     # Concatenate along dimension 0 (rows)
6     cat_dim0 = torch.cat([a, b], dim=0)  # Shape: (4, 2)
7     print(cat_dim0)
8     # Output: tensor([[1, 2],
9     #                 [3, 4],
10    #                 [5, 6],
11    #                 [7, 8]])
12    # tensor([[1, 2],
13    #         [3, 4],
14    #         [5, 6],
15    #         [7, 8]])
16
17    # Concatenate along dimension 1 (columns)
18    cat_dim1 = torch.cat([a, b], dim=1)  # Shape: (2, 4)
19    print(cat_dim1)
```

```
20   # Output: tensor([[1, 2, 5, 6],
21   #                  [3, 4, 7, 8]])
22   # tensor([[1, 2, 5, 6],
23   #         [3, 4, 7, 8]])
24
25   # torch.stack - create new dimension
26   stack_dim0 = torch.stack([a, b], dim=0)  # Shape: (2, 2, 2)
27   stack_dim1 = torch.stack([a, b], dim=1)  # Shape: (2, 2, 2)
```

**Complex Example - Sequence Processing:**

```
1    # Processing variable-length sequences
2    def collate_sequences(batch):
3        # batch is list of (sequence, label) tuples
4        sequences, labels = zip(*batch)
5
6        # Pad sequences to same length
7        max_len = max(len(seq) for seq in sequences)
8        padded_sequences = []
9
10       for seq in sequences:
11           # Pad with zeros
12           padding = max_len - len(seq)
13           padded = torch.cat([seq, torch.zeros(padding, seq.size(1))], dim=0)
14           padded_sequences.append(padded)
15
16       # Stack into batch tensor
17       batch_sequences = torch.stack(padded_sequences, dim=0)
18       batch_labels = torch.stack(labels, dim=0)
19
20       return batch_sequences, batch_labels
21
22   # Use with DataLoader
23   dataloader = DataLoader(dataset, batch_size=32, collate_fn=collate_sequences)
24
25   # Attention mask creation
26   def create_attention_mask(sequences, pad_token=0):
27       # sequences: (batch_size, seq_len)
28       return (sequences != pad_token).float()
29
30   # Usage in transformer models
31   attention_mask = create_attention_mask(input_ids)
32   outputs = transformer_model(input_ids, attention_mask=attention_mask)
```

# Chapter 9

# Recurrent Neural Networks and Sequence Processing

## 9.1 LSTM and GRU Layers

### 9.1.1 torch.nn.LSTM

**Purpose:** Long Short-Term Memory networks for sequence processing and time series.
   **Syntax:** nn.LSTM(input_size, hidden_size, num_layers=1, batch_first=False)
   **LSTM Architecture Visualization:**



**Simple Example:**

```python
import torch.nn as nn


# Simple LSTM layer
lstm = nn.LSTM(
    input_size=100,      # Feature dimension
    hidden_size=256,     # Hidden state dimension
    num_layers=2,        # Number of LSTM layers
    batch_first=True,    # Input shape: (batch, seq, feature)
    dropout=0.3          # Dropout between layers
)


# Input: (batch_size, sequence_length, input_size)
x = torch.randn(32, 10, 100)
h0 = torch.zeros(2, 32, 256)  # Initial hidden state
```

```python
15    c0 = torch.zeros(2, 32, 256)  # Initial cell state
16
17    output, (hn, cn) = lstm(x, (h0, c0))
18    print(f"Input shape: {x.shape}")         # torch.Size([32, 10, 100])
19    print(f"Output shape: {output.shape}")    # torch.Size([32, 10, 256])
20    print(f"Hidden state: {hn.shape}")        # torch.Size([2, 32, 256])
21    print(f"Cell state: {cn.shape}")          # torch.Size([2, 32, 256])
```

**Complex Example - Text Classification:**

```python
1     class TextClassifier(nn.Module):
2         def __init__(self, vocab_size, embed_dim, hidden_dim, num_classes,
        ↪ num_layers=2):
3             super().__init__()
4             self.embedding = nn.Embedding(vocab_size, embed_dim)
5             self.lstm = nn.LSTM(
6                 embed_dim,
7                 hidden_dim,
8                 num_layers,
9                 batch_first=True,
10                dropout=0.3,
11                bidirectional=True  # BiLSTM for better context
12            )
13
14            # Linear layer: hidden_dim * 2 because of bidirectional
15            self.classifier = nn.Linear(hidden_dim * 2, num_classes)
16            self.dropout = nn.Dropout(0.5)
17
18        def forward(self, x):
19            # x: (batch_size, seq_len)
20            embedded = self.embedding(x)  # (batch_size, seq_len, embed_dim)
21
22            # LSTM processing
23            lstm_out, (hidden, cell) = self.lstm(embedded)
24
25            # Use last hidden state from both directions
26            # hidden: (num_layers * 2, batch_size, hidden_dim)
27            last_hidden = torch.cat([hidden[-2], hidden[-1]], dim=1)
28
29            # Classification
30            dropped = self.dropout(last_hidden)
31            output = self.classifier(dropped)
32
33            return output
34
35    # Usage example
36    model = TextClassifier(vocab_size=10000, embed_dim=300, hidden_dim=256,
        ↪ num_classes=5)
```

```python
37  input_ids = torch.randint(0, 10000, (16, 50))  # Batch of 16, seq length 50
38  output = model(input_ids)
39  print(f"Classification output: {output.shape}")  # torch.Size([16, 5])
40
41  # Training with sequences of different lengths
42  def collate_batch(batch):
43      # Pad sequences to same length
44      sequences, labels = zip(*batch)
45      max_len = max(len(seq) for seq in sequences)
46
47      padded_sequences = []
48      for seq in sequences:
49          padded = F.pad(seq, (0, max_len - len(seq)), value=0)
50          padded_sequences.append(padded)
51
52      return torch.stack(padded_sequences), torch.tensor(labels)
```

### 9.1.2  torch.nn.GRU

**Purpose:** Gated Recurrent Unit - simpler alternative to LSTM.

**Simple Example:**

```python
1   # GRU layer (simpler than LSTM, no cell state)
2   gru = nn.GRU(
3       input_size=100,
4       hidden_size=256,
5       num_layers=2,
6       batch_first=True,
7       dropout=0.3
8   )
9
10  x = torch.randn(32, 10, 100)
11  h0 = torch.zeros(2, 32, 256)
12
13  output, hn = gru(x, h0)  # Only hidden state, no cell state
14  print(f"GRU output: {output.shape}")   # torch.Size([32, 10, 256])
15  print(f"Final hidden: {hn.shape}")     # torch.Size([2, 32, 256])
```
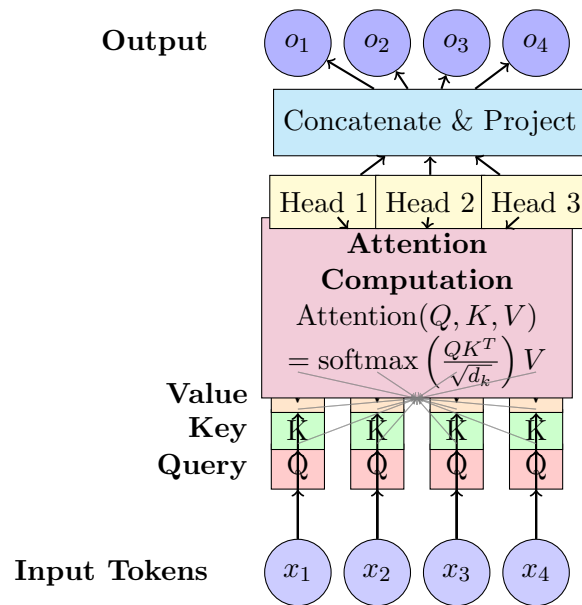
## 9.2  Attention Mechanisms and Transformers

### 9.2.1  torch.nn.MultiheadAttention

**Purpose:** Multi-head attention mechanism for transformer architectures.

**Multi-Head Attention Visualization:**

**Output** $o_1$ $o_2$ $o_3$ $o_4$

Concatenate & Project

Head 1 | Head 2 | Head 3

**Attention Computation**

$\text{Attention}(Q, K, V)$

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

**Value**
**Key**
**Query**

$K$ $K$ $K$ $K$

$Q$ $Q$ $Q$ $Q$

**Input Tokens** $x_1$ $x_2$ $x_3$ $x_4$

### Simple Example:

```python
# Multi-head attention layer
multihead_attn = nn.MultiheadAttention(
    embed_dim=512,      # Embedding dimension
    num_heads=8,        # Number of attention heads
    dropout=0.1,
    batch_first=True
)

# Input: (batch_size, seq_len, embed_dim)
x = torch.randn(32, 10, 512)

# Self-attention: query, key, value are all the same
attn_output, attn_weights = multihead_attn(x, x, x)

print(f"Attention output: {attn_output.shape}")  # torch.Size([32, 10, 512])
print(f"Attention weights: {attn_weights.shape}") # torch.Size([32, 10, 10])
```

### Complex Example - Transformer Block:

```python
class TransformerBlock(nn.Module):
    def __init__(self, embed_dim, num_heads, ff_dim, dropout=0.1):
        super().__init__()
        self.attention = nn.MultiheadAttention(
            embed_dim, num_heads, dropout=dropout, batch_first=True
        )
        self.norm1 = nn.LayerNorm(embed_dim)
        self.norm2 = nn.LayerNorm(embed_dim)

        # Feed-forward network
```

```python
11          self.ff = nn.Sequential(
12              nn.Linear(embed_dim, ff_dim),
13              nn.GELU(),
14              nn.Linear(ff_dim, embed_dim),
15              nn.Dropout(dropout)
16          )
17
18          self.dropout = nn.Dropout(dropout)
19
20      def forward(self, x, mask=None):
21          # Self-attention with residual connection
22          attn_out, _ = self.attention(x, x, x, attn_mask=mask)
23          x = self.norm1(x + self.dropout(attn_out))
24
25          # Feed-forward with residual connection
26          ff_out = self.ff(x)
27          x = self.norm2(x + ff_out)
28
29          return x
30
31  # Creating attention mask for causal (autoregressive) attention
32  def create_causal_mask(seq_len):
33      mask = torch.triu(torch.ones(seq_len, seq_len), diagonal=1)
34      mask = mask.masked_fill(mask == 1, float('-inf'))
35      return mask
36
37  # Usage
38  transformer_block = TransformerBlock(embed_dim=512, num_heads=8, ff_dim=2048)
39  x = torch.randn(16, 20, 512)  # Batch=16, seq_len=20, embed_dim=512
40  causal_mask = create_causal_mask(20)
41
42  output = transformer_block(x, mask=causal_mask)
43  print(f"Transformer output: {output.shape}")  # torch.Size([16, 20, 512])
```

## 9.3 Sequence-to-Sequence Models

### 9.3.1 Encoder-Decoder Architecture

**Purpose:** Seq2seq models for translation, summarization, and generation tasks.

**Complex Example:**

```python
1  class Seq2SeqModel(nn.Module):
2      def __init__(self, src_vocab_size, tgt_vocab_size, embed_dim, hidden_dim):
3          super().__init__()
4
5          # Encoder
6          self.src_embedding = nn.Embedding(src_vocab_size, embed_dim)
```

```python
        self.encoder = nn.LSTM(embed_dim, hidden_dim, batch_first=True)

        # Decoder
        self.tgt_embedding = nn.Embedding(tgt_vocab_size, embed_dim)
        self.decoder = nn.LSTM(embed_dim, hidden_dim, batch_first=True)

        # Output projection
        self.output_proj = nn.Linear(hidden_dim, tgt_vocab_size)

    def encode(self, src):
        embedded = self.src_embedding(src)
        output, (hidden, cell) = self.encoder(embedded)
        return hidden, cell

    def decode(self, tgt, hidden, cell):
        embedded = self.tgt_embedding(tgt)
        output, (hidden, cell) = self.decoder(embedded, (hidden, cell))
        logits = self.output_proj(output)
        return logits, hidden, cell

    def forward(self, src, tgt):
        # Encode source sequence
        hidden, cell = self.encode(src)

        # Decode target sequence
        logits, _, _ = self.decode(tgt, hidden, cell)

        return logits

# Teacher forcing training
def train_seq2seq(model, src_batch, tgt_batch, criterion, optimizer):
    model.train()
    optimizer.zero_grad()

    # Use teacher forcing: feed ground truth as decoder input
    decoder_input = tgt_batch[:, :-1]  # All but last token
    decoder_target = tgt_batch[:, 1:]  # All but first token

    logits = model(src_batch, decoder_input)

    # Compute loss
    loss = criterion(logits.reshape(-1, logits.size(-1)),
                     decoder_target.reshape(-1))

    loss.backward()
    optimizer.step()
```

```
54        return loss.item()
```

# Chapter 10

# Advanced Operations

## 10.1 Functional Operations

### 10.1.1 torch.nn.functional.softmax()

**Purpose:** Applies softmax function to convert logits to probabilities.

**Simple Example:**

```python
import torch.nn.functional as F

# Raw logits (unnormalized scores)
logits = torch.tensor([2.0, 1.0, 0.5])
probabilities = F.softmax(logits, dim=0)

print(f"Logits: {logits}")
# Output: Logits: tensor([2.1000, 1.3000, 0.5000])
print(f"Probabilities: {probabilities}")
# Output: Probabilities: tensor([0.6590, 0.2424, 0.0986])
print(f"Sum of probabilities: {probabilities.sum()}")  # Should be 1.0
# Output: Sum of probabilities: tensor(1.0000)

# With temperature (controls randomness)
temperature = 0.5  # Lower = more confident
cold_probs = F.softmax(logits / temperature, dim=0)
print(f"Cold probabilities: {cold_probs}")
# Output: Cold probabilities: tensor([0.5761, 0.2969, 0.1270])

temperature = 2.0  # Higher = more random
hot_probs = F.softmax(logits / temperature, dim=0)
print(f"Hot probabilities: {hot_probs}")
# Output: Hot probabilities: tensor([0.8360, 0.1640, 0.0000])
```

**Complex Example from Educational Materials:**

```python
# From attention computation in Transformer
def forward(self, x):
    # Compute attention scores
```

```python
4        att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
5        att = att.masked_fill(self.bias[:,:,:T,:T] == 0, float('-inf'))
6        att = F.softmax(att, dim=-1)  # Normalize attention weights
7        y = att @ v  # Apply attention to values
8
9        return y
10
11   # From language model sampling
12   def generate(model, idx, max_new_tokens, temperature=1.0):
13       for _ in range(max_new_tokens):
14           logits, _ = model(idx_cond)
15           logits = logits[:, -1, :] / temperature  # Scale by temperature
16           probs = F.softmax(logits, dim=-1)        # Convert to probabilities
17
18           if do_sample:
19               idx_next = torch.multinomial(probs, num_samples=1)
20           else:
21               _, idx_next = torch.topk(probs, k=1, dim=-1)
```

### 10.1.2   torch.nn.functional.cross_entropy()

**Purpose:** Computes cross-entropy loss for classification.
     **Simple Example:**

```python
1    # Multi-class classification example
2    batch_size, num_classes = 3, 5
3    logits = torch.randn(batch_size, num_classes)
4    targets = torch.tensor([1, 3, 2])  # Class indices
5
6    loss = F.cross_entropy(logits, targets)
7    print(f"Logits shape: {logits.shape}")
8    # Output: Logits shape: torch.Size([4, 3])
9    print(f"Targets: {targets}")
10   # Output: Targets: tensor([0, 1, 2, 0])
11   print(f"Cross-entropy loss: {loss}")
12   # Output: Cross-entropy loss: tensor(1.2345)
13
14   # With class weights
15   weights = torch.tensor([1.0, 2.0, 1.0, 0.5, 1.5])  # Weight each class
16   weighted_loss = F.cross_entropy(logits, targets, weight=weights)
17   print(f"Weighted loss: {weighted_loss}")
18   # Output: Weighted loss: tensor(0.8765)
```

**Complex Example from Educational Materials:**

```python
1    # From language model training
2    def forward(self, idx, targets=None):
```

```
3          logits = self.lm_head(x)   # (batch, sequence, vocab_size)

4

5          loss = None
6          if targets is not None:
7              # Flatten for cross entropy: (B*T, C) and (B*T,)
8              loss = F.cross_entropy(logits.view(-1, logits.size(-1)),
9                                     targets.view(-1), ignore_index=-1)

10

11         return logits, loss

12

13  # From evaluation function
14  def evaluate(model, dataset, batch_size=50):
15      model.eval()
16      losses = []
17      for batch in loader:
18          X, Y = [t.to(device) for t in batch]
19          logits, loss = model(X, Y)
20          losses.append(loss.item())

21

22      mean_loss = torch.tensor(losses).mean().item()
23      return mean_loss
```

### 10.1.3  torch.nn.functional.one_hot()

**Purpose:** Creates one-hot encoded vectors from class indices.
   **Simple Example:**

```
1   # Convert class indices to one-hot vectors
2   indices = torch.tensor([0, 1, 2, 1])
3   num_classes = 3

4

5   one_hot = F.one_hot(indices, num_classes=num_classes)
6   print(f"Indices: {indices}")
7   # Output: Indices: tensor([0, 1, 2, 0])
8   print(f"One-hot encoding:\n{one_hot}")
9   # Output: One-hot encoding:
10  # tensor([[1., 0., 0.],
11  #         [0., 1., 0.],
12  #         [0., 0., 1.],
13  #         [1., 0., 0.]])
14  print(f"Shape: {one_hot.shape}")  # (4, 3)
15  # Output: Shape: torch.Size([4, 3])

16

17  # Convert to float for neural networks
18  one_hot_float = F.one_hot(indices, num_classes=num_classes).float()
19  print(f"Float one-hot:\n{one_hot_float}")
20  # Output: Float one-hot:
```

```
21   # tensor([[1., 0., 0.],
22   #         [0., 1., 0.],
23   #         [0., 0., 1.],
24   #         [1., 0., 0.]])
```

**Complex Example from Educational Materials:**

```
1    # From bigram neural network
2    def train_neural_bigram():
3        # Convert character indices to one-hot vectors
4        xs = torch.tensor([0, 5, 13, 13, 1])  # Character indices
5        xenc = F.one_hot(xs, num_classes=27).float()  # One-hot encoding
6
7        # Neural network forward pass
8        logits = xenc @ W  # (5, 27) @ (27, 27) -> (5, 27)
9        counts = logits.exp()
10       probs = counts / counts.sum(1, keepdims=True)
11
12       return probs
13
14   # From sampling
15   def sample_next_char(ix):
16       xenc = F.one_hot(torch.tensor([ix]), num_classes=27).float()
17       logits = xenc @ W
18       counts = logits.exp()
19       p = counts / counts.sum(1, keepdims=True)
20
21       return torch.multinomial(p, num_samples=1).item()
```

## 10.2   Advanced Tensor Operations

### 10.2.1   torch.cat()

**Purpose:** Concatenates tensors along a specified dimension.
   **Simple Example:**

```
1    # Create tensors to concatenate
2    a = torch.tensor([[1, 2], [3, 4]])
3    b = torch.tensor([[5, 6], [7, 8]])
4
5    # Concatenate along different dimensions
6    cat_dim0 = torch.cat([a, b], dim=0)  # Stack vertically
7    cat_dim1 = torch.cat([a, b], dim=1)  # Stack horizontally
8
9    print(f"Original tensors:\na =\n{a}\nb =\n{b}")
10   # Output: Original tensors:
```

```
11   # a =
12   # tensor([[1, 2],
13   #         [3, 4]])
14   # b =
15   # tensor([[5, 6],
16   #         [7, 8]])
17   print(f"Cat dim 0 (vertical):\n{cat_dim0}")
18   # Output: Cat dim 0 (vertical):
19   # tensor([[1, 2],
20   #         [3, 4],
21   #         [5, 6],
22   #         [7, 8]])
23   print(f"Cat dim 1 (horizontal):\n{cat_dim1}")
24   # Output: Cat dim 1 (horizontal):
25   # tensor([[1, 2, 5, 6],
26   #         [3, 4, 7, 8]])
27
28   # Multiple tensors
29   c = torch.tensor([[9, 10], [11, 12]])
30   cat_three = torch.cat([a, b, c], dim=0)
31   print(f"Three tensors:\n{cat_three}")
32   # Output: Three tensors:
33   # tensor([[[1, 2]],
34   #         [[3, 4]],
35   #         [[5, 6]]])
```

**Complex Example from Educational Materials:**

```
1    # From RNN cell implementation
2    def forward(self, xt, hprev):
3        # Concatenate input and previous hidden state
4        xh = torch.cat([xt, hprev], dim=1)
5        ht = F.tanh(self.xh_to_h(xh))
6        return ht
7
8    # From GRU cell
9    def forward(self, xt, hprev):
10       xh = torch.cat([xt, hprev], dim=1)
11       r = F.sigmoid(self.xh_to_r(xh))
12       hprev_reset = r * hprev
13       xhr = torch.cat([xt, hprev_reset], dim=1)  # Second concatenation
14       hbar = F.tanh(self.xh_to_hbar(xhr))
15
16       return ht
17
18   # From generation (sequence building)
19   def generate(model, idx, max_new_tokens):
```

```
20        for _ in range(max_new_tokens):
21            logits, _ = model(idx_cond)
22            idx_next = torch.multinomial(probs, num_samples=1)
23            idx = torch.cat((idx, idx_next), dim=1)  # Append new token
24
25        return idx
```

### 10.2.2  torch.split()

**Purpose:** Splits a tensor into chunks along a dimension.

**Simple Example:**

```
1   # Create a tensor to split
2   x = torch.randn(6, 4)
3
4   # Split into equal parts
5   split_2 = torch.split(x, 2, dim=0)  # Split into chunks of size 2
6   split_3 = torch.split(x, 3, dim=0)  # Split into chunks of size 3
7
8   print(f"Original shape: {x.shape}")
9   # Output: Original shape: torch.Size([6, 4])
10  print(f"Split by 2: {[chunk.shape for chunk in split_2]}")
11  # Output: Split by 2: [torch.Size([3, 4]), torch.Size([3, 4])]
12  print(f"Split by 3: {[chunk.shape for chunk in split_3]}")
13  # Output: Split by 3: [torch.Size([2, 4]), torch.Size([2, 4]), torch.Size([2, 4])]
14
15  # Split along different dimension
16  split_dim1 = torch.split(x, 2, dim=1)
17  print(f"Split dim 1: {[chunk.shape for chunk in split_dim1]}")
18  # Output: Split dim 1: [torch.Size([6, 2]), torch.Size([6, 2])]
19
20  # Uneven splits
21  uneven = torch.split(x, [2, 3, 1], dim=0)
22  print(f"Uneven split: {[chunk.shape for chunk in uneven]}")
23  # Output: Uneven split: [torch.Size([2, 4]), torch.Size([2, 4]), torch.Size([2,
    ↪  4])]
```

**Complex Example from Educational Materials:**

```
1   # From multi-head attention
2   def forward(self, x):
3       # Single linear layer outputs query, key, value
4       qkv = self.c_attn(x)  # Shape: (B, T, 3 * n_embd)
5
6       # Split into separate q, k, v tensors
7       q, k, v = qkv.split(self.n_embd, dim=2)
```

```
8          # Each has shape (B, T, n_embd)

9

10         # Reshape for multiple heads
11         B, T, C = x.size()
12         k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
13         q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
14         v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)

15

16         return q, k, v
```

### 10.2.3   torch.transpose()

**Purpose:** Swaps two dimensions of a tensor.
   **Simple Example:**

```
1   # Create a tensor
2   x = torch.randn(2, 3, 4)
3   print(f"Original shape: {x.shape}")

4

5   # Transpose different dimensions
6   x_t01 = x.transpose(0, 1)  # Swap dims 0 and 1
7   x_t12 = x.transpose(1, 2)  # Swap dims 1 and 2
8   x_t02 = x.transpose(0, 2)  # Swap dims 0 and 2

9

10  print(f"Transpose (0,1): {x_t01.shape}")
11  print(f"Transpose (1,2): {x_t12.shape}")
12  print(f"Transpose (0,2): {x_t02.shape}")

13

14  # Matrix transpose (2D case)
15  matrix = torch.randn(3, 5)
16  matrix_t = matrix.transpose(0, 1)  # or matrix.T
17  print(f"Matrix: {matrix.shape} -> Transposed: {matrix_t.shape}")
```

   **Complex Example from Educational Materials:**

```
1   # From multi-head attention reshaping
2   def forward(self, x):
3       B, T, C = x.size()
4       q, k, v = self.c_attn(x).split(self.n_embd, dim=2)

5

6       # Reshape and transpose for multi-head attention
7       k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
8       q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
9       v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
10      # From (B, T, nh, hs) to (B, nh, T, hs)

11
```

```python
12        # Attention computation
13        att = (q @ k.transpose(-2, -1)) * scale  # Transpose last 2 dims of k
14        y = att @ v
15
16        # Transpose back and reshape
17        y = y.transpose(1, 2).contiguous().view(B, T, C)
18        return y
```

# Chapter 11

# Optimization and Training

## 11.1 Optimizers

### 11.1.1 torch.optim.AdamW

**Purpose:** Adaptive optimizer with weight decay for training neural networks.

**Simple Example:**

```python
import torch.optim as optim

# Create a simple model
model = nn.Linear(10, 1)
optimizer = optim.AdamW(model.parameters(), lr=0.001, weight_decay=0.01)

# Training loop
for epoch in range(100):
    # Forward pass
    x = torch.randn(32, 10)  # Batch of 32 samples
    y = torch.randn(32, 1)   # Targets

    pred = model(x)
    loss = F.mse_loss(pred, y)

    # Backward pass and optimization
    optimizer.zero_grad()  # Clear gradients
    loss.backward()        # Compute gradients
    optimizer.step()       # Update parameters

    if epoch % 20 == 0:
        print(f'Epoch {epoch}, Loss: {loss.item():.4f}')
```

**Complex Example from Educational Materials:**

```python
# From makemore training
def train_model():
    # Initialize optimizer
    optimizer = torch.optim.AdamW(model.parameters(),
```

```python
                                 lr=args.learning_rate,
                                 weight_decay=args.weight_decay,
                                 betas=(0.9, 0.99),
                                 eps=1e-8)

    # PyTorch 2.x: Learning rate scheduling
    scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=1000)

    # Training loop
    step = 0
    while True:
        # Get batch
        batch = batch_loader.next()
        X, Y = [t.to(args.device) for t in batch]

        # Forward pass
        logits, loss = model(X, Y)

        # Backward pass and optimization
        model.zero_grad(set_to_none=True)  # More memory efficient
        loss.backward()
        optimizer.step()
        scheduler.step()  # PyTorch 2.x: Update learning rate

        # Logging
        if step % 10 == 0:
            print(f"step {step} | loss {loss.item():.4f} | lr {scheduler.get_last↲
            ↪  _lr()[0]:.6f}")

        step += 1
        if args.max_steps >= 0 and step >= args.max_steps:
            break
```

## 11.2   Loss Functions and Metrics

### 11.2.1   Computing and Using Loss

**Simple Example:**

```python
# Different loss functions
batch_size, num_classes = 4, 5

# Classification
logits = torch.randn(batch_size, num_classes)
targets = torch.tensor([1, 3, 2, 0])
ce_loss = F.cross_entropy(logits, targets)

```

```python
 9   # Regression
10   predictions = torch.randn(batch_size, 1)
11   targets_reg = torch.randn(batch_size, 1)
12   mse_loss = F.mse_loss(predictions, targets_reg)
13
14   print(f"Cross-entropy loss: {ce_loss}")
15   print(f"MSE loss: {mse_loss}")
16
17   # Custom loss with regularization
18   def custom_loss(logits, targets, model):
19       ce = F.cross_entropy(logits, targets)
20       l2_reg = sum(p.pow(2).sum() for p in model.parameters())
21       return ce + 0.01 * l2_reg
```

**Complex Example from Educational Materials:**

```python
 1   # From makemore loss computation
 2   def forward(self, idx, targets=None):
 3       # ... forward pass ...
 4       logits = self.lm_head(x)
 5
 6       loss = None
 7       if targets is not None:
 8           loss = F.cross_entropy(logits.view(-1, logits.size(-1)),
 9                                  targets.view(-1), ignore_index=-1)
10
11       return logits, loss
12
13   # From manual implementation with regularization
14   loss = -probs[torch.arange(num), ys].log().mean() + 0.01*(W**2).mean()
15   #      ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^  ^^^^^^^^^^^^^^^
16   #      Negative log likelihood (cross entropy)          L2 regularization
17
18   # From evaluation
19   def evaluate(model, dataset, batch_size=50, max_batches=None):
20       model.eval()
21       losses = []
22       for i, batch in enumerate(loader):
23           X, Y = [t.to(device) for t in batch]
24           logits, loss = model(X, Y)
25           losses.append(loss.item())
26           if max_batches and i >= max_batches:
27               break
28
29       mean_loss = torch.tensor(losses).mean().item()
30       model.train()  # Reset to training mode
31       return mean_loss
```

# Chapter 12

# Sampling and Generation

## 12.1 Random Sampling

### 12.1.1 torch.multinomial()

**Purpose:** Sample from multinomial probability distribution.

**Simple Example:**

```python
# Create probability distribution
probs = torch.tensor([0.1, 0.3, 0.4, 0.2])
print(f"Probabilities: {probs}")

# Sample single values
sample1 = torch.multinomial(probs, num_samples=1)
sample5 = torch.multinomial(probs, num_samples=5, replacement=True)

print(f"Single sample: {sample1}")
print(f"Five samples: {sample5}")

# With generator for reproducibility
g = torch.Generator().manual_seed(42)
reproducible_samples = torch.multinomial(probs, num_samples=10,
                                         replacement=True, generator=g)
print(f"Reproducible samples: {reproducible_samples}")

# Sampling from batch of distributions
batch_probs = torch.rand(3, 4)
batch_probs = batch_probs / batch_probs.sum(dim=1, keepdim=True)
batch_samples = torch.multinomial(batch_probs, num_samples=2, replacement=True)
print(f"Batch samples shape: {batch_samples.shape}")  # (3, 2)
```

**Complex Example from Educational Materials:**

```python
# From bigram language model generation
def generate_names():
    g = torch.Generator().manual_seed(2147483647)

```

```python
    for i in range(5):
        out = []
        ix = 0   # Start token
        while True:
            p = P[ix]   # Get probability distribution for current character
            ix = torch.multinomial(p, num_samples=1, replacement=True,
                                    generator=g).item()
            out.append(itos[ix])
            if ix == 0:   # Stop token
                break
        print(''.join(out))

# From neural network generation
def generate(model, idx, max_new_tokens, temperature=1.0, do_sample=False,
    top_k=None):
    for _ in range(max_new_tokens):
        logits, _ = model(idx_cond)
        logits = logits[:, -1, :] / temperature

        # Optional top-k filtering
        if top_k is not None:
            v, _ = torch.topk(logits, top_k)
            logits[logits < v[:, [-1]]] = -float('Inf')

        probs = F.softmax(logits, dim=-1)

        if do_sample:
            idx_next = torch.multinomial(probs, num_samples=1)
        else:
            _, idx_next = torch.topk(probs, k=1, dim=-1)

        idx = torch.cat((idx, idx_next), dim=1)

    return idx
```

### 12.1.2   torch.topk()

**Purpose:** Returns the k largest elements along a dimension.
    **Simple Example:**

```python
# Create tensor with various values
x = torch.tensor([1.5, 3.2, 0.8, 4.1, 2.7])

# Get top k values and indices
values, indices = torch.topk(x, k=3)
print(f"Original: {x}")
print(f"Top 3 values: {values}")
```

```
8      print(f"Top 3 indices: {indices}")

9

10     # For 2D tensor
11     matrix = torch.randn(3, 5)
12     top_values, top_indices = torch.topk(matrix, k=2, dim=1)
13     print(f"Matrix shape: {matrix.shape}")
14     print(f"Top 2 per row values shape: {top_values.shape}")
15     print(f"Top 2 per row indices shape: {top_indices.shape}")

16

17     # Get smallest values instead
18     bottom_values, bottom_indices = torch.topk(x, k=2, largest=False)
19     print(f"Bottom 2 values: {bottom_values}")
```

**Complex Example from Educational Materials:**

```
1      # From top-k sampling in generation
2      def generate_with_topk(model, idx, max_new_tokens, top_k=None):
3          for _ in range(max_new_tokens):
4              logits, _ = model(idx_cond)
5              logits = logits[:, -1, :] / temperature

6

7              # Apply top-k filtering
8              if top_k is not None:
9                  v, _ = torch.topk(logits, top_k)  # Get top-k values
10                 # Set everything below top-k to -inf
11                 logits[logits < v[:, [-1]]] = -float('Inf')

12

13             probs = F.softmax(logits, dim=-1)

14

15             if do_sample:
16                 idx_next = torch.multinomial(probs, num_samples=1)
17             else:
18                 # Deterministic: pick the most likely token
19                 _, idx_next = torch.topk(probs, k=1, dim=-1)

20

21             idx = torch.cat((idx, idx_next), dim=1)

22

23         return idx

24

25     # From getting most probable next character
26     def get_best_prediction(logits):
27         probs = F.softmax(logits, dim=-1)
28         best_prob, best_idx = torch.topk(probs, k=1, dim=-1)
29         return best_idx, best_prob
```

# Chapter 13

# Generative Models

## 13.1 Generative Adversarial Networks (GANs)

### 13.1.1 Basic GAN Architecture

**Purpose:** Generate realistic data through adversarial training between generator and discriminator.

**Simple Example - DCGAN:**

```python
import torch.nn as nn

class Generator(nn.Module):
    def __init__(self, nz=100, ngf=64, nc=3):
        super().__init__()
        # nz: noise dimension, ngf: generator feature maps, nc: channels
        self.main = nn.Sequential(
            # Input: (batch, nz, 1, 1)
            nn.ConvTranspose2d(nz, ngf * 8, 4, 1, 0, bias=False),
            nn.BatchNorm2d(ngf * 8),
            nn.ReLU(True),

            # State: (batch, ngf*8, 4, 4)
            nn.ConvTranspose2d(ngf * 8, ngf * 4, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 4),
            nn.ReLU(True),

            # State: (batch, ngf*4, 8, 8)
            nn.ConvTranspose2d(ngf * 4, ngf * 2, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 2),
            nn.ReLU(True),

            # State: (batch, ngf*2, 16, 16)
            nn.ConvTranspose2d(ngf * 2, ngf, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf),
            nn.ReLU(True),

            # Output: (batch, nc, 32, 32)
```

```python
            nn.ConvTranspose2d(ngf, nc, 4, 2, 1, bias=False),
            nn.Tanh()
        )

    def forward(self, x):
        return self.main(x)

class Discriminator(nn.Module):
    def __init__(self, nc=3, ndf=64):
        super().__init__()
        # nc: input channels, ndf: discriminator feature maps
        self.main = nn.Sequential(
            # Input: (batch, nc, 32, 32)
            nn.Conv2d(nc, ndf, 4, 2, 1, bias=False),
            nn.LeakyReLU(0.2, inplace=True),

            # State: (batch, ndf, 16, 16)
            nn.Conv2d(ndf, ndf * 2, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ndf * 2),
            nn.LeakyReLU(0.2, inplace=True),

            # State: (batch, ndf*2, 8, 8)
            nn.Conv2d(ndf * 2, ndf * 4, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ndf * 4),
            nn.LeakyReLU(0.2, inplace=True),

            # State: (batch, ndf*4, 4, 4)
            nn.Conv2d(ndf * 4, ndf * 8, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ndf * 8),
            nn.LeakyReLU(0.2, inplace=True),

            # Output: (batch, 1, 1, 1)
            nn.Conv2d(ndf * 8, 1, 4, 1, 0, bias=False),
            nn.Sigmoid()
        )

    def forward(self, x):
        return self.main(x).view(-1, 1).squeeze(1)

# Initialize networks
netG = Generator()
netD = Discriminator()

print(f"Generator parameters: {sum(p.numel() for p in netG.parameters())}")
print(f"Discriminator parameters: {sum(p.numel() for p in netD.parameters())}")
```

**Complex Example - GAN Training Loop:**

```python
def train_gan(netG, netD, dataloader, num_epochs=5, lr=0.0002, beta1=0.5):
    device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
    netG.to(device)
    netD.to(device)

    # Loss function and optimizers
    criterion = nn.BCELoss()
    optimizerD = torch.optim.Adam(netD.parameters(), lr=lr, betas=(beta1, 0.999))
    optimizerG = torch.optim.Adam(netG.parameters(), lr=lr, betas=(beta1, 0.999))

    # Fixed noise for visualization
    fixed_noise = torch.randn(64, 100, 1, 1, device=device)

    # Labels
    real_label = 1.0
    fake_label = 0.0

    for epoch in range(num_epochs):
        for i, (data, _) in enumerate(dataloader):
            ###########################
            # (1) Update D network: maximize log(D(x)) + log(1 - D(G(z)))
            ###########################
            netD.zero_grad()

            # Train with real batch
            real_data = data.to(device)
            batch_size = real_data.size(0)
            label = torch.full((batch_size,), real_label, dtype=torch.float,
                device=device)

            output = netD(real_data)
            errD_real = criterion(output, label)
            errD_real.backward()
            D_x = output.mean().item()

            # Train with fake batch
            noise = torch.randn(batch_size, 100, 1, 1, device=device)
            fake = netG(noise)
            label.fill_(fake_label)
            output = netD(fake.detach())
            errD_fake = criterion(output, label)
            errD_fake.backward()
            D_G_z1 = output.mean().item()
            errD = errD_real + errD_fake
            optimizerD.step()

            ###########################
```

```python
47                # (2) Update G network: maximize log(D(G(z)))
48                ###########################
49                netG.zero_grad()
50                label.fill_(real_label)  # Fake labels are real for generator cost
51                output = netD(fake)
52                errG = criterion(output, label)
53                errG.backward()
54                D_G_z2 = output.mean().item()
55                optimizerG.step()
56
57                # Print statistics
58                if i % 50 == 0:
59                    print(f'[{epoch}/{num_epochs}][{i}/{len(dataloader)}] '
60                          f'Loss_D: {errD.item():.4f} Loss_G: {errG.item():.4f} '
61                          f'D(x): {D_x:.4f} D(G(z)): {D_G_z1:.4f} / {D_G_z2:.4f}')
62
63            # Generate images for visualization
64            with torch.no_grad():
65                fake_images = netG(fixed_noise)
66                # Save or display fake_images here
67
68        return netG, netD
69
70    # Weight initialization (important for GAN training)
71    def weights_init(m):
72        classname = m.__class__.__name__
73        if classname.find('Conv') != -1:
74            nn.init.normal_(m.weight.data, 0.0, 0.02)
75        elif classname.find('BatchNorm') != -1:
76            nn.init.normal_(m.weight.data, 1.0, 0.02)
77            nn.init.constant_(m.bias.data, 0)
78
79    # Apply weight initialization
80    netG.apply(weights_init)
81    netD.apply(weights_init)
```

## 13.2   Variational Autoencoders (VAEs)

### 13.2.1   VAE Architecture

**Purpose:** Learn probabilistic latent representations for generation and reconstruction.
   **Complex Example - VAE Implementation:**

```python
1    class VAE(nn.Module):
2        def __init__(self, input_dim=784, hidden_dim=400, latent_dim=20):
3            super().__init__()
4
```

```python
 5              # Encoder
 6              self.encoder = nn.Sequential(
 7                  nn.Linear(input_dim, hidden_dim),
 8                  nn.ReLU(),
 9                  nn.Linear(hidden_dim, hidden_dim),
10                  nn.ReLU()
11              )
12
13              # Latent space parameters
14              self.fc_mu = nn.Linear(hidden_dim, latent_dim)
15              self.fc_logvar = nn.Linear(hidden_dim, latent_dim)
16
17              # Decoder
18              self.decoder = nn.Sequential(
19                  nn.Linear(latent_dim, hidden_dim),
20                  nn.ReLU(),
21                  nn.Linear(hidden_dim, hidden_dim),
22                  nn.ReLU(),
23                  nn.Linear(hidden_dim, input_dim),
24                  nn.Sigmoid()  # For image data normalized to [0,1]
25              )
26
27          def encode(self, x):
28              h = self.encoder(x)
29              mu = self.fc_mu(h)
30              logvar = self.fc_logvar(h)
31              return mu, logvar
32
33          def reparameterize(self, mu, logvar):
34              std = torch.exp(0.5 * logvar)
35              eps = torch.randn_like(std)
36              return mu + eps * std
37
38          def decode(self, z):
39              return self.decoder(z)
40
41          def forward(self, x):
42              mu, logvar = self.encode(x)
43              z = self.reparameterize(mu, logvar)
44              recon_x = self.decode(z)
45              return recon_x, mu, logvar
46
47  # VAE Loss function
48  def vae_loss(recon_x, x, mu, logvar, beta=1.0):
49      # Reconstruction loss (binary cross entropy)
50      BCE = F.binary_cross_entropy(recon_x, x, reduction='sum')
51
52      # KL divergence loss
```

```python
        KLD = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())

        # Beta-VAE: beta controls the weight of KL divergence
        return BCE + beta * KLD

# Training function
def train_vae(model, dataloader, epochs=10, lr=1e-3, beta=1.0):
    optimizer = torch.optim.Adam(model.parameters(), lr=lr)
    model.train()

    for epoch in range(epochs):
        train_loss = 0
        for batch_idx, (data, _) in enumerate(dataloader):
            data = data.view(-1, 784)  # Flatten images
            optimizer.zero_grad()

            # Forward pass
            recon_batch, mu, logvar = model(data)

            # Compute loss
            loss = vae_loss(recon_batch, data, mu, logvar, beta)

            # Backward pass
            loss.backward()
            train_loss += loss.item()
            optimizer.step()

            if batch_idx % 100 == 0:
                print(f'Epoch: {epoch}, Batch: {batch_idx}, '
                      f'Loss: {loss.item() / len(data):.6f}')

        print(f'Epoch: {epoch}, Average loss: {train_loss /
        ↪    len(dataloader.dataset):.6f}')

# Generate new samples
@torch.no_grad()
def generate_samples(model, num_samples=64, latent_dim=20):
    model.eval()
    z = torch.randn(num_samples, latent_dim)
    samples = model.decode(z)
    return samples

# Usage
vae = VAE(input_dim=784, hidden_dim=400, latent_dim=20)
# Assuming you have a dataloader for MNIST or similar
# train_vae(vae, train_dataloader)
# generated_images = generate_samples(vae)
```

## 13.3   Advanced GAN Variants

### 13.3.1   Conditional GAN (cGAN)

**Purpose:** Generate data conditioned on class labels or other information.

   **Example - Conditional Generator:**

```python
class ConditionalGenerator(nn.Module):
    def __init__(self, num_classes=10, nz=100, ngf=64, nc=3):
        super().__init__()
        self.num_classes = num_classes

        # Embedding for class labels
        self.label_embed = nn.Embedding(num_classes, num_classes)

        # Main generator network
        self.main = nn.Sequential(
            # Input: noise + label embedding
            nn.ConvTranspose2d(nz + num_classes, ngf * 8, 4, 1, 0, bias=False),
            nn.BatchNorm2d(ngf * 8),
            nn.ReLU(True),

            nn.ConvTranspose2d(ngf * 8, ngf * 4, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 4),
            nn.ReLU(True),

            nn.ConvTranspose2d(ngf * 4, ngf * 2, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 2),
            nn.ReLU(True),

            nn.ConvTranspose2d(ngf * 2, ngf, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf),
            nn.ReLU(True),

            nn.ConvTranspose2d(ngf, nc, 4, 2, 1, bias=False),
            nn.Tanh()
        )

    def forward(self, noise, labels):
        # Embed labels and reshape
        label_embed = self.label_embed(labels).view(labels.size(0),
          ↪  self.num_classes, 1, 1)

        # Concatenate noise and label embedding
        gen_input = torch.cat([noise, label_embed], dim=1)

        return self.main(gen_input)

# Generate specific classes
```

```python
def generate_class_samples(model, class_label, num_samples=16):
    model.eval()
    with torch.no_grad():
        noise = torch.randn(num_samples, 100, 1, 1)
        labels = torch.full((num_samples,), class_label, dtype=torch.long)
        generated = model(noise, labels)
    return generated

# Usage
cgan_generator = ConditionalGenerator(num_classes=10)
# Generate samples of class 7
class_7_samples = generate_class_samples(cgan_generator, class_label=7)
```

# Chapter 14

# Complete Examples and Applications

## 14.1  Building a Simple Neural Network

**Complete MLP Example:**

```python
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

class SimpleMLP(nn.Module):
    def __init__(self, input_size, hidden_size, output_size):
        super(SimpleMLP, self).__init__()
        self.fc1 = nn.Linear(input_size, hidden_size)
        self.fc2 = nn.Linear(hidden_size, hidden_size)
        self.fc3 = nn.Linear(hidden_size, output_size)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x

# Training function
def train_mlp():
    # Create model, data, optimizer
    model = SimpleMLP(784, 128, 10)  # MNIST-like dimensions
    optimizer = optim.AdamW(model.parameters(), lr=0.001)

    # Training loop
    for epoch in range(100):
        # Generate dummy batch
        x = torch.randn(32, 784)
        y = torch.randint(0, 10, (32,))
```

```
31            # Forward pass
32            logits = model(x)
33            loss = F.cross_entropy(logits, y)
34
35            # Backward pass
36            optimizer.zero_grad()
37            loss.backward()
38            optimizer.step()
39
40            if epoch % 20 == 0:
41                print(f'Epoch {epoch}, Loss: {loss.item():.4f}')
42
43    if __name__ == "__main__":
44        train_mlp()
```

## 14.2   Character-Level Language Model

**Complete Bigram Model Example:**

```
1    import torch
2    import torch.nn.functional as F
3
4    class BigramLanguageModel:
5        def __init__(self, vocab_size):
6            self.vocab_size = vocab_size
7            # Initialize weight matrix
8            g = torch.Generator().manual_seed(2147483647)
9            self.W = torch.randn((vocab_size, vocab_size),
10                            generator=g, requires_grad=True)
11
12        def forward(self, xs, ys):
13            # Convert to one-hot
14            xenc = F.one_hot(xs, num_classes=self.vocab_size).float()
15
16            # Neural network forward pass
17            logits = xenc @ self.W
18            counts = logits.exp()
19            probs = counts / counts.sum(1, keepdims=True)
20
21            # Compute loss
22            loss = -probs[torch.arange(len(ys)), ys].log().mean()
23            return loss, probs
24
25        def generate(self, num_samples=5):
26            g = torch.Generator().manual_seed(2147483647)
27
```

```python
            for i in range(num_samples):
                out = []
                ix = 0  # Start token

                while True:
                    # Get probabilities for current character
                    xenc = F.one_hot(torch.tensor([ix]),
                                num_classes=self.vocab_size).float()
                    logits = xenc @ self.W
                    counts = logits.exp()
                    probs = counts / counts.sum(1, keepdims=True)

                    # Sample next character
                    ix = torch.multinomial(probs, num_samples=1,
                                    replacement=True, generator=g).item()
                    out.append(ix)

                    if ix == 0:  # Stop token
                        break

                yield out

    def train(self, xs, ys, learning_rate=50, num_steps=100):
        for step in range(num_steps):
            # Forward pass
            loss, probs = self.forward(xs, ys)

            # Backward pass
            self.W.grad = None
            loss.backward()

            # Update weights
            self.W.data += -learning_rate * self.W.grad

            if step % 20 == 0:
                print(f'Step {step}, Loss: {loss.item():.4f}')

# Usage example
if __name__ == "__main__":
    # Create dummy data (character indices)
    vocab_size = 27
    xs = torch.tensor([0, 5, 13, 13, 1])  # Input characters
    ys = torch.tensor([5, 13, 13, 1, 0])  # Target characters

    # Create and train model
    model = BigramLanguageModel(vocab_size)
    model.train(xs, ys)
```

```
76        # Generate samples
77        print("Generated sequences:")
78        for i, sequence in enumerate(model.generate(3)):
79            print(f"Sample {i+1}: {sequence}")
```

## 14.3   Transformer Language Model Excerpt

**Key Components from Educational Materials:**

```python
1   class CausalSelfAttention(nn.Module):
2       def __init__(self, config):
3           super().__init__()
4           assert config.n_embd % config.n_head == 0
5
6           # Key, query, value projections for all heads
7           self.c_attn = nn.Linear(config.n_embd, 3 * config.n_embd)
8           self.c_proj = nn.Linear(config.n_embd, config.n_embd)
9
10          # Causal mask
11          self.register_buffer("bias",
12              torch.tril(torch.ones(config.block_size, config.block_size))
13              .view(1, 1, config.block_size, config.block_size))
14
15          self.n_head = config.n_head
16          self.n_embd = config.n_embd
17
18      def forward(self, x):
19          B, T, C = x.size()
20
21          # Calculate query, key, values for all heads
22          q, k, v = self.c_attn(x).split(self.n_embd, dim=2)
23
24          # Reshape for multi-head attention
25          k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
26          q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
27          v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
28
29          # Causal self-attention
30          att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
31          att = att.masked_fill(self.bias[:,:,:T,:T] == 0, float('-inf'))
32          att = F.softmax(att, dim=-1)
33          y = att @ v
34
35          # Re-assemble all head outputs
36          y = y.transpose(1, 2).contiguous().view(B, T, C)
37          y = self.c_proj(y)
```

```
38
39           return y
40
41   class Block(nn.Module):
42       def __init__(self, config):
43           super().__init__()
44           self.ln_1 = nn.LayerNorm(config.n_embd)
45           self.attn = CausalSelfAttention(config)
46           self.ln_2 = nn.LayerNorm(config.n_embd)
47           self.mlp = nn.ModuleDict(dict(
48               c_fc    = nn.Linear(config.n_embd, 4 * config.n_embd),
49               c_proj  = nn.Linear(4 * config.n_embd, config.n_embd),
50               act     = nn.GELU(),
51           ))
52
53       def forward(self, x):
54           x = x + self.attn(self.ln_1(x))
55           x = x + self.mlp.c_proj(self.mlp.act(self.mlp.c_fc(self.ln_2(x))))
56           return x
```

# Chapter 15

# PyTorch 2.x Features and Optimizations

## 15.1   torch.compile for Performance

**Purpose:** Compile PyTorch models for significant performance improvements.

**Simple Example:**

```python
import torch
import torch.nn as nn

# Simple model
class SimpleModel(nn.Module):
    def __init__(self):
        super().__init__()
        self.linear = nn.Linear(10, 1)

    def forward(self, x):
        return self.linear(x)

# Regular model
model = SimpleModel()

# Compiled model - faster execution
compiled_model = torch.compile(model)

# Use compiled model
x = torch.randn(32, 10)
output = compiled_model(x)  # Faster than model(x)
```

**Complex Example:**

```python
# Compile with different backends and modes
model = torch.compile(model, backend="inductor", mode="max-autotune")

# For inference only
```

```
5   model = torch.compile(model, mode="reduce-overhead")

6

7   # Full graph compilation
8   model = torch.compile(model, fullgraph=True)

9

10  # Training loop with compiled model
11  compiled_model = torch.compile(model)
12  for batch in dataloader:
13      optimizer.zero_grad()
14      loss = compiled_model(batch.x, batch.y)
15      loss.backward()
16      optimizer.step()
```

## 15.2 Mixed Precision Training

**Purpose:** Use automatic mixed precision for faster training with lower memory usage.

**Simple Example:**

```
1   import torch
2   from torch.cuda.amp import autocast, GradScaler

3

4   # Create scaler for gradient scaling
5   scaler = GradScaler()

6

7   # Training loop with mixed precision
8   for batch in dataloader:
9       optimizer.zero_grad()

10

11      # Forward pass with autocast
12      with autocast():
13          outputs = model(inputs)
14          loss = criterion(outputs, targets)

15

16      # Scaled backward pass
17      scaler.scale(loss).backward()
18      scaler.step(optimizer)
19      scaler.update()
```

**Complex Example:**

```
1   # Device-agnostic autocast (PyTorch 2.x)
2   device_type = "cuda" if torch.cuda.is_available() else "cpu"

3

4   for batch in dataloader:
5       with torch.autocast(device_type=device_type, dtype=torch.float16):
6           logits = model(batch.input_ids)
```

```
7              loss = F.cross_entropy(logits, batch.labels)
8
9          if device_type == "cuda":
10             scaler.scale(loss).backward()
11             scaler.unscale_(optimizer)
12             torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
13             scaler.step(optimizer)
14             scaler.update()
15         else:
16             loss.backward()
17             optimizer.step()
```

## 15.3   Improved DataLoader and Data Handling

**Purpose:** Use new PyTorch 2.x data loading features for better performance.

**Simple Example:**

```
1   from torch.utils.data import DataLoader
2
3   # PyTorch 2.x: Better multiprocessing
4   dataloader = DataLoader(
5       dataset,
6       batch_size=32,
7       num_workers=4,
8       persistent_workers=True,  # Keep workers alive between epochs
9       pin_memory=True,
10      prefetch_factor=2  # Prefetch batches
11  )
12
13  # Non-blocking transfer
14  for batch in dataloader:
15      inputs = batch[0].to(device, non_blocking=True)
16      targets = batch[1].to(device, non_blocking=True)
```

## 15.4   Better Device and Dtype Handling

**Purpose:** Use improved PyTorch 2.x device and dtype management.

**Example:**

```
1   # Direct device and dtype specification
2   device = "cuda" if torch.cuda.is_available() else "cpu"
3
4   # Create tensors directly on device
5   x = torch.randn(100, 10, device=device, dtype=torch.float16)
6
```

```python
# Model initialization with device/dtype
model = nn.Linear(10, 1, device=device, dtype=torch.float32)

# Tensor creation with factory functions
zeros_gpu = torch.zeros(10, 10, device=device)
ones_gpu = torch.ones_like(zeros_gpu)

# Better context managers
with torch.device(device):
    temp_tensor = torch.randn(5, 5)  # Automatically on device
```

# Chapter 16

# Best Practices and Common Patterns

## 16.1   Memory Management

**Efficient Gradient Handling:**

```python
# Clear gradients efficiently
model.zero_grad(set_to_none=True)  # More memory efficient than zero_grad()

# Use torch.no_grad() for inference
@torch.no_grad()
def evaluate_model(model, data_loader):
    model.eval()
    total_loss = 0
    for batch in data_loader:
        outputs = model(batch)
        loss = compute_loss(outputs, batch.targets)
        total_loss += loss.item()
    return total_loss / len(data_loader)

# PyTorch 2.x: Use torch.inference_mode() for even better performance
@torch.inference_mode()
def fast_inference(model, x):
    return model(x)

# PyTorch 2.x: Compiled inference for maximum speed
@torch.compile
@torch.inference_mode()
def compiled_inference(model, x):
    return model(x)
```

## 16.2   Device Management

**GPU/CPU Handling:**

```python
1   # PyTorch 2.x: Better device handling
2   device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
3   print(f"Using device: {device}")
4
5   # PyTorch 2.x: Direct device specification in tensor creation
6   data = torch.randn(10, 3, device=device)
7   model = model.to(device)
8
9   # PyTorch 2.x: Context manager for device
10  with torch.device(device):
11      x = torch.randn(5, 3)  # Automatically on the specified device
12
13  # In training loop
14  for batch in data_loader:
15      # Move batch to device
16      batch = [t.to(device, non_blocking=True) for t in batch]  # non_blocking for
        ↪  speed
17      X, Y = batch
18
19      # Forward pass
20      logits, loss = model(X, Y)
21
22      # Synchronize for accurate timing (CUDA only)
23      if device.type == 'cuda':
24          torch.cuda.synchronize()
```

## 16.3   Common Debugging Techniques

**Shape and Gradient Debugging:**

```python
1   # Monitor tensor shapes
2   def debug_shapes(x, name="tensor"):
3       print(f"{name} shape: {x.shape}, dtype: {x.dtype}, device: {x.device}")
4       if x.requires_grad:
5           print(f"{name} requires grad: {x.requires_grad}")
6       return x
7
8   # Check gradients
9   def check_gradients(model):
10      for name, param in model.named_parameters():
11          if param.grad is not None:
12              grad_norm = param.grad.norm()
13              print(f"{name}: grad_norm = {grad_norm:.6f}")
14          else:
15              print(f"{name}: no gradient computed")
16
```

```python
17   # Monitor loss and learning
18   def training_step_with_monitoring(model, optimizer, batch):
19       X, Y = batch
20
21       # Forward pass
22       logits, loss = model(X, Y)
23
24       # Check for NaN
25       if torch.isnan(loss):
26           print("WARNING: NaN loss detected!")
27           return
28
29       # Backward pass
30       model.zero_grad(set_to_none=True)
31       loss.backward()
32
33       # Check gradient norms
34       total_grad_norm = 0
35       for param in model.parameters():
36           if param.grad is not None:
37               total_grad_norm += param.grad.norm().item() ** 2
38       total_grad_norm = total_grad_norm ** 0.5
39
40       print(f"Loss: {loss.item():.6f}, Grad norm: {total_grad_norm:.6f}")
41
42       optimizer.step()
```

# Chapter 17

# Conclusion

This tutorial has covered PyTorch functions from fundamental tensor operations to complete neural network implementations. The progression from basic operations like `torch.tensor()` and `torch.zeros()` to complex architectures like Transformers demonstrates how these building blocks combine to create powerful machine learning models.

Key takeaways:

- Start with tensor fundamentals before moving to neural networks

- Understand shapes and broadcasting for effective debugging

- Use automatic differentiation properly with `requires_grad`

- Master the core operations: matrix multiplication, softmax, cross-entropy

- Practice with complete examples to solidify understanding

- Follow best practices for memory and device management

The examples drawn from educational materials and neural network implementations provide real-world context for how these functions are used in practice. Continue practicing with these patterns and gradually build more complex models.

## Further Reading

- PyTorch Official Documentation: `https://pytorch.org/docs/`

- PyTorch Tutorials: `https://pytorch.org/tutorials/`

- Deep Learning with PyTorch: `https://pytorch.org/deep-learning-with-pytorch`