

Deep Learning & Neural Networks

A Comprehensive Guide to Modern Deep Learning

Contents

1	Introduction to Deep Learning	4
1.1	Why Deep Learning?	4
2	Artificial Neural Networks (ANNs)	4
2.1	The Perceptron	4
2.1.1	Mathematical Formulation	4
2.2	Activation Functions	4
2.2.1	Common Activation Functions	4
2.2.2	Python Implementation	5
2.3	Multi-Layer Perceptron (MLP)	5
2.3.1	Forward Propagation	5
2.3.2	Loss Functions	6
2.4	Backpropagation	6
2.4.1	Gradient Computation	6
2.4.2	Python Implementation from Scratch	6
2.4.3	Using TensorFlow/Keras	9
3	Optimization Algorithms	10
3.1	Gradient Descent Variants	10
3.1.1	Batch Gradient Descent	10
3.1.2	Stochastic Gradient Descent (SGD)	10
3.1.3	Mini-Batch Gradient Descent	10
3.2	Advanced Optimizers	10
3.2.1	Momentum	10
3.2.2	RMSprop	10
3.2.3	Adam (Adaptive Moment Estimation)	11
4	Regularization Techniques	11
4.1	L2 Regularization (Weight Decay)	11
4.2	Dropout	12
4.3	Batch Normalization	12
5	Convolutional Neural Networks (CNNs)	12
5.1	Convolutional Layer	12
5.1.1	Convolution Operation	12
5.1.2	Properties of Convolution	13
5.2	Pooling Layer	13
5.3	CNN Architecture	13
5.3.1	Python Implementation	13
5.4	Famous CNN Architectures	14
5.4.1	LeNet-5 (1998)	14

5.4.2	AlexNet (2012)	14
5.4.3	VGGNet (2014)	15
5.4.4	ResNet (2015)	15
6	Recurrent Neural Networks (RNNs)	16
6.1	Basic RNN	16
6.1.1	Mathematical Formulation	16
6.1.2	Backpropagation Through Time (BPTT)	16
6.2	Long Short-Term Memory (LSTM)	16
6.2.1	LSTM Equations	16
6.3	Gated Recurrent Unit (GRU)	17
6.3.1	GRU Equations	17
6.3.2	Python Implementation	17
6.4	Bidirectional RNNs	18
7	Attention Mechanisms	18
7.1	Attention Score	19
7.2	Multi-Head Attention	19
8	Transformers	20
8.1	Transformer Architecture	20
9	Autoencoders	21
9.1	Vanilla Autoencoder	21
9.2	Variational Autoencoder (VAE)	21
10	Generative Adversarial Networks (GANs)	24
10.1	GAN Objective	24
11	Transfer Learning	26
11.1	Approaches	26
12	Best Practices & Tips	27
12.1	Training Strategies	27
12.2	Model Evaluation	28
12.3	Hyperparameter Tuning	29
12.4	Debugging Neural Networks	30
13	Advanced Topics	30
13.1	Neural Architecture Search (NAS)	30
13.2	Meta-Learning	31
13.3	Neural ODEs	31
13.4	Graph Neural Networks (GNNs)	31
14	Practical Projects	31
14.1	Image Classification	31
14.2	Text Generation with RNN	32
15	Resources & Further Reading	34
15.1	Books	34
15.2	Online Courses	34
15.3	Research Papers	34

1 Introduction to Deep Learning

Deep learning is a subset of machine learning based on artificial neural networks with multiple layers. These networks can learn hierarchical representations of data, making them particularly effective for complex tasks like image recognition, natural language processing, and game playing.

1.1 Why Deep Learning?

- **Feature Learning:** Automatically learns features from raw data
- **Scalability:** Performance improves with more data
- **End-to-End Learning:** Can optimize entire pipelines jointly
- **Transfer Learning:** Pre-trained models can be adapted to new tasks

2 Artificial Neural Networks (ANNs)

2.1 The Perceptron

The perceptron is the simplest neural network unit, inspired by biological neurons.

2.1.1 Mathematical Formulation

Given input vector $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ and weights $\mathbf{w} = [w_1, w_2, \dots, w_n]^T$:

$$z = \mathbf{w}^T \mathbf{x} + b = \sum_{i=1}^n w_i x_i + b \quad (1)$$

$$y = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases} \quad (2)$$

2.2 Activation Functions

Activation functions introduce non-linearity, enabling networks to learn complex patterns.

2.2.1 Common Activation Functions

1. Sigmoid:

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad \sigma'(z) = \sigma(z)(1 - \sigma(z)) \quad (3)$$

2. Hyperbolic Tangent (\tanh):

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}, \quad \tanh'(z) = 1 - \tanh^2(z) \quad (4)$$

3. ReLU (Rectified Linear Unit):

$$\text{ReLU}(z) = \max(0, z), \quad \text{ReLU}'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases} \quad (5)$$

4. Leaky ReLU:

$$\text{LeakyReLU}(z) = \max(\alpha z, z), \quad \alpha \approx 0.01 \quad (6)$$

5. Softmax (for output layer):

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad (7)$$

2.2.2 Python Implementation

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def sigmoid(z):
5     return 1 / (1 + np.exp(-z))
6
7 def tanh(z):
8     return np.tanh(z)
9
10 def relu(z):
11     return np.maximum(0, z)
12
13 def leaky_relu(z, alpha=0.01):
14     return np.where(z > 0, z, alpha * z)
15
16 # Visualization
17 z = np.linspace(-5, 5, 100)
18 plt.figure(figsize=(12, 8))
19
20 plt.subplot(2, 2, 1)
21 plt.plot(z, sigmoid(z))
22 plt.title('Sigmoid')
23 plt.grid(True)
24
25 plt.subplot(2, 2, 2)
26 plt.plot(z, tanh(z))
27 plt.title('Tanh')
28 plt.grid(True)
29
30 plt.subplot(2, 2, 3)
31 plt.plot(z, relu(z))
32 plt.title('ReLU')
33 plt.grid(True)
34
35 plt.subplot(2, 2, 4)
36 plt.plot(z, leaky_relu(z))
37 plt.title('Leaky ReLU')
38 plt.grid(True)
39
40 plt.tight_layout()
41 plt.show()
```

2.3 Multi-Layer Perceptron (MLP)

An MLP consists of an input layer, one or more hidden layers, and an output layer.

2.3.1 Forward Propagation

For a network with L layers:

$$\mathbf{z}^{[l]} = \mathbf{W}^{[l]} \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]} \quad (8)$$

$$\mathbf{a}^{[l]} = g^{[l]}(\mathbf{z}^{[l]}) \quad (9)$$

where:

- $\mathbf{z}^{[l]}$ is the weighted sum at layer l
- $\mathbf{a}^{[l]}$ is the activation at layer l
- $g^{[l]}$ is the activation function
- $\mathbf{W}^{[l]}$ is the weight matrix
- $\mathbf{b}^{[l]}$ is the bias vector

2.3.2 Loss Functions

Mean Squared Error (Regression):

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (10)$$

Cross-Entropy Loss (Classification):

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log(\hat{y}_{ik}) \quad (11)$$

2.4 Backpropagation

Backpropagation computes gradients using the chain rule.

2.4.1 Gradient Computation

For layer l :

$$\delta^{[l]} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \odot g'^{[l]}(\mathbf{z}^{[l]}) \quad (12)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]} (\mathbf{a}^{[l-1]})^T \quad (13)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]} \quad (14)$$

Backpropagation recursion:

$$\delta^{[l]} = (\mathbf{W}^{[l+1]})^T \delta^{[l+1]} \odot g'^{[l]}(\mathbf{z}^{[l]}) \quad (15)$$

where \odot denotes element-wise multiplication.

2.4.2 Python Implementation from Scratch

```

1 import numpy as np
2
3 class NeuralNetwork:
4     def __init__(self, layer_dims):
5         """
6             layer_dims: list containing dimensions of each layer
7             Example: [784, 128, 64, 10] for MNIST
8         """
9         self.parameters = {}
10        self.L = len(layer_dims) - 1 # number of layers
11

```

```

12     # Initialize weights and biases
13     for l in range(1, self.L + 1):
14         self.parameters[f'W{l}'] = np.random.randn(
15             layer_dims[l], layer_dims[l-1]
16         ) * 0.01
17         self.parameters[f'b{l}'] = np.zeros((layer_dims[l], 1))
18
19     def relu(self, Z):
20         return np.maximum(0, Z)
21
22     def relu_derivative(self, Z):
23         return (Z > 0).astype(float)
24
25     def softmax(self, Z):
26         exp_Z = np.exp(Z - np.max(Z, axis=0, keepdims=True))
27         return exp_Z / np.sum(exp_Z, axis=0, keepdims=True)
28
29     def forward_propagation(self, X):
30         cache = {'A0': X}
31         A = X
32
33         # Hidden layers with ReLU
34         for l in range(1, self.L):
35             W = self.parameters[f'W{l}']
36             b = self.parameters[f'b{l}']
37             Z = np.dot(W, A) + b
38             A = self.relu(Z)
39             cache[f'Z{l}'] = Z
40             cache[f'A{l}'] = A
41
42         # Output layer with softmax
43         W = self.parameters[f'W{self.L}']
44         b = self.parameters[f'b{self.L}']
45         Z = np.dot(W, A) + b
46         A = self.softmax(Z)
47         cache[f'Z{self.L}'] = Z
48         cache[f'A{self.L}'] = A
49
50         return A, cache
51
52     def compute_cost(self, AL, Y):
53         m = Y.shape[1]
54         cost = -np.sum(Y * np.log(AL + 1e-8)) / m
55         return cost
56
57     def backward_propagation(self, cache, Y):
58         gradients = {}
59         m = Y.shape[1]
60
61         # Output layer gradient
62         dZ = cache[f'A{self.L}'] - Y
63         gradients[f'dW{self.L}'] = np.dot(dZ, cache[f'A{self.L-1}'].T) /
64             m
65         gradients[f'db{self.L}'] = np.sum(dZ, axis=1, keepdims=True) /
66             m
67
68         # Hidden layers gradients
69         for l in reversed(range(1, self.L)):

```

```

68     dA = np.dot(self.parameters[f'W{1+1}'].T, dZ)
69     dZ = dA * self.relu_derivative(cache[f'Z{1}'])
70     gradients[f'dW{1}'] = np.dot(dZ, cache[f'A{1-1}'].T) / m
71     gradients[f'db{1}'] = np.sum(dZ, axis=1, keepdims=True) / m
72
73     return gradients
74
75 def update_parameters(self, gradients, learning_rate):
76     for l in range(1, self.L + 1):
77         self.parameters[f'W{l}'] -= learning_rate * gradients[f'dW{l}']
78         self.parameters[f'b{l}'] -= learning_rate * gradients[f'db{l}']
79
80 def train(self, X, Y, epochs, learning_rate):
81     costs = []
82     for epoch in range(epochs):
83         # Forward propagation
84         AL, cache = self.forward_propagation(X)
85
86         # Compute cost
87         cost = self.compute_cost(AL, Y)
88         costs.append(cost)
89
90         # Backward propagation
91         gradients = self.backward_propagation(cache, Y)
92
93         # Update parameters
94         self.update_parameters(gradients, learning_rate)
95
96         if epoch % 100 == 0:
97             print(f"Epoch {epoch}, Cost: {cost:.4f}")
98
99     return costs
100
101 def predict(self, X):
102     AL, _ = self.forward_propagation(X)
103     return np.argmax(AL, axis=0)
104
105 # Example usage with MNIST-like data
106 # Generate dummy data
107 np.random.seed(42)
108 X_train = np.random.randn(784, 1000) # 1000 samples, 784 features
109 Y_train = np.zeros((10, 1000)) # 10 classes
110 Y_train[np.random.randint(0, 10, 1000), np.arange(1000)] = 1
111
112 # Create and train network
113 nn = NeuralNetwork([784, 128, 64, 10])
114 costs = nn.train(X_train, Y_train, epochs=1000, learning_rate=0.01)
115
116 # Plot learning curve
117 plt.plot(costs)
118 plt.xlabel('Epoch')
119 plt.ylabel('Cost')
120 plt.title('Learning Curve')
121 plt.grid(True)
122 plt.show()

```

2.4.3 Using TensorFlow/Keras

```
1 import tensorflow as tf
2 from tensorflow import keras
3 from tensorflow.keras import layers
4 from sklearn.datasets import load_digits
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7
8 # Load data
9 digits = load_digits()
10 X, y = digits.data, digits.target
11
12 # Preprocess
13 scaler = StandardScaler()
14 X_scaled = scaler.fit_transform(X)
15 X_train, X_test, y_train, y_test = train_test_split(
16     X_scaled, y, test_size=0.2, random_state=42
17 )
18
19 # Build model
20 model = keras.Sequential([
21     layers.Dense(128, activation='relu', input_shape=(64,)),
22     layers.Dropout(0.3),
23     layers.Dense(64, activation='relu'),
24     layers.Dropout(0.3),
25     layers.Dense(10, activation='softmax')
26 ])
27
28 # Compile
29 model.compile(
30     optimizer='adam',
31     loss='sparse_categorical_crossentropy',
32     metrics=['accuracy']
33 )
34
35 # Train
36 history = model.fit(
37     X_train, y_train,
38     epochs=50,
39     batch_size=32,
40     validation_split=0.2,
41     verbose=1
42 )
43
44 # Evaluate
45 test_loss, test_acc = model.evaluate(X_test, y_test, verbose=0)
46 print(f"\nTest accuracy: {test_acc:.4f}")
47
48 # Plot training history
49 plt.figure(figsize=(12, 4))
50 plt.subplot(1, 2, 1)
51 plt.plot(history.history['loss'], label='Training Loss')
52 plt.plot(history.history['val_loss'], label='Validation Loss')
53 plt.xlabel('Epoch')
54 plt.ylabel('Loss')
55 plt.legend()
56 plt.grid(True)
```

```

57 plt.subplot(1, 2, 2)
58 plt.plot(history.history['accuracy'], label='Training Accuracy')
59 plt.plot(history.history['val_accuracy'], label='Validation Accuracy')
60 plt.xlabel('Epoch')
61 plt.ylabel('Accuracy')
62 plt.legend()
63 plt.grid(True)
64 plt.tight_layout()
65 plt.show()

```

3 Optimization Algorithms

3.1 Gradient Descent Variants

3.1.1 Batch Gradient Descent

Updates parameters using the entire dataset:

$$\mathbf{W} := \mathbf{W} - \alpha \nabla_{\mathbf{W}} \mathcal{L} \quad (16)$$

3.1.2 Stochastic Gradient Descent (SGD)

Updates parameters using one sample at a time:

$$\mathbf{W} := \mathbf{W} - \alpha \nabla_{\mathbf{W}} \mathcal{L}^{(i)} \quad (17)$$

3.1.3 Mini-Batch Gradient Descent

Updates using small batches of size m :

$$\mathbf{W} := \mathbf{W} - \alpha \frac{1}{m} \sum_{i=1}^m \nabla_{\mathbf{W}} \mathcal{L}^{(i)} \quad (18)$$

3.2 Advanced Optimizers

3.2.1 Momentum

Accelerates gradient descent by accumulating velocity:

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \nabla_{\mathbf{W}} \mathcal{L} \quad (19)$$

$$\mathbf{W} := \mathbf{W} - \alpha \mathbf{v}_t \quad (20)$$

Typical: $\beta = 0.9$

3.2.2 RMSprop

Adapts learning rate based on gradient magnitude:

$$\mathbf{s}_t = \beta \mathbf{s}_{t-1} + (1 - \beta) (\nabla_{\mathbf{W}} \mathcal{L})^2 \quad (21)$$

$$\mathbf{W} := \mathbf{W} - \frac{\alpha}{\sqrt{\mathbf{s}_t + \epsilon}} \nabla_{\mathbf{W}} \mathcal{L} \quad (22)$$

3.2.3 Adam (Adaptive Moment Estimation)

Combines momentum and RMSprop:

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla_{\mathbf{W}} \mathcal{L} \quad (23)$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla_{\mathbf{W}} \mathcal{L})^2 \quad (24)$$

$$\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t}, \quad \hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_2^t} \quad (25)$$

$$\mathbf{W} := \mathbf{W} - \frac{\alpha}{\sqrt{\hat{\mathbf{v}}_t} + \epsilon} \hat{\mathbf{m}}_t \quad (26)$$

Typical: $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$

```

1 # Comparing optimizers in Keras
2 optimizers_list = [
3     ('SGD', keras.optimizers.SGD(learning_rate=0.01)),
4     ('Momentum', keras.optimizers.SGD(learning_rate=0.01, momentum=0.9)),
5     ('RMSprop', keras.optimizers.RMSprop(learning_rate=0.001)),
6     ('Adam', keras.optimizers.Adam(learning_rate=0.001))
7 ]
8
9 histories = {}
10 for name, optimizer in optimizers_list:
11     model = keras.Sequential([
12         layers.Dense(128, activation='relu', input_shape=(64,)),
13         layers.Dense(64, activation='relu'),
14         layers.Dense(10, activation='softmax')
15     ])
16
17     model.compile(optimizer=optimizer,
18                   loss='sparse_categorical_crossentropy',
19                   metrics=['accuracy'])
20
21     history = model.fit(X_train, y_train, epochs=30,
22                          batch_size=32, validation_split=0.2,
23                          verbose=0)
24     histories[name] = history
25
26 # Plot comparison
27 plt.figure(figsize=(12, 4))
28 for name, history in histories.items():
29     plt.plot(history.history['val_accuracy'], label=name)
30 plt.xlabel('Epoch')
31 plt.ylabel('Validation Accuracy')
32 plt.title('Optimizer Comparison')
33 plt.legend()
34 plt.grid(True)
35 plt.show()
```

4 Regularization Techniques

4.1 L2 Regularization (Weight Decay)

Add penalty to loss function:

$$\mathcal{L}_{\text{reg}} = \mathcal{L} + \frac{\lambda}{2} \sum_{l=1}^L \|\mathbf{W}^{[l]}\|_F^2 \quad (27)$$

4.2 Dropout

Randomly deactivate neurons during training with probability p :

$$\mathbf{a}^{[l]} = \mathbf{a}^{[l]} \odot \mathbf{d}^{[l]} \quad (28)$$

where $\mathbf{d}^{[l]} \sim \text{Bernoulli}(1 - p)$

During inference: $\mathbf{a}^{[l]} = (1 - p)\mathbf{a}^{[l]}$

4.3 Batch Normalization

Normalize activations within mini-batch:

$$\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m z_i \quad (29)$$

$$\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (z_i - \mu_{\mathcal{B}})^2 \quad (30)$$

$$\hat{z}_i = \frac{z_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad (31)$$

$$\tilde{z}_i = \gamma \hat{z}_i + \beta \quad (32)$$

where γ and β are learnable parameters.

```

1 # Regularization in Keras
2 model = keras.Sequential([
3     layers.Dense(128, activation='relu', input_shape=(64,)),
4         kernel_regularizer=keras.regularizers.l2(0.01)),
5     layers.BatchNormalization(),
6     layers.Dropout(0.5),
7
8     layers.Dense(64, activation='relu',
9                 kernel_regularizer=keras.regularizers.l2(0.01)),
10    layers.BatchNormalization(),
11    layers.Dropout(0.3),
12
13    layers.Dense(10, activation='softmax')
14])

```

5 Convolutional Neural Networks (CNNs)

CNNs are specialized for processing grid-like data, particularly images.

5.1 Convolutional Layer

5.1.1 Convolution Operation

For input \mathbf{I} and kernel \mathbf{K} :

$$(\mathbf{I} * \mathbf{K})(i, j) = \sum_m \sum_n \mathbf{I}(i + m, j + n) \mathbf{K}(m, n) \quad (33)$$

Output dimensions:

$$\text{Output size} = \left\lfloor \frac{n + 2p - f}{s} \right\rfloor + 1 \quad (34)$$

where:

- n : input size
- p : padding
- f : filter size
- s : stride

5.1.2 Properties of Convolution

- **Parameter Sharing**: Same filter applied across entire input
- **Translation Invariance**: Detects features regardless of position
- **Local Connectivity**: Each neuron connects to local region

5.2 Pooling Layer

Reduces spatial dimensions while retaining important features.

Max Pooling:

$$y_{i,j} = \max_{m,n \in \mathcal{R}} x_{i+m,j+n} \quad (35)$$

Average Pooling:

$$y_{i,j} = \frac{1}{|\mathcal{R}|} \sum_{m,n \in \mathcal{R}} x_{i+m,j+n} \quad (36)$$

5.3 CNN Architecture

Typical architecture:

$$\text{INPUT} \rightarrow [\text{CONV} \rightarrow \text{RELU} \rightarrow \text{POOL}]^* \rightarrow \text{FC} \rightarrow \text{SOFTMAX} \quad (37)$$

5.3.1 Python Implementation

```

1 import tensorflow as tf
2 from tensorflow.keras import layers, models
3 from tensorflow.keras.datasets import mnist
4
5 # Load MNIST dataset
6 (X_train, y_train), (X_test, y_test) = mnist.load_data()
7 X_train = X_train.reshape(-1, 28, 28, 1).astype('float32') / 255.0
8 X_test = X_test.reshape(-1, 28, 28, 1).astype('float32') / 255.0
9
10 # Build CNN
11 model = models.Sequential([
12     # First convolutional block
13     layers.Conv2D(32, (3, 3), activation='relu',
14                 input_shape=(28, 28, 1), padding='same'),
15     layers.BatchNormalization(),
16     layers.MaxPooling2D((2, 2)),
17
18     # Second convolutional block
19     layers.Conv2D(64, (3, 3), activation='relu', padding='same'),
20     layers.BatchNormalization(),
21     layers.MaxPooling2D((2, 2)),
22
23     # Third convolutional block

```

```

24     layers.Conv2D(128, (3, 3), activation='relu', padding='same'),
25     layers.BatchNormalization(),
26     layers.MaxPooling2D((2, 2)),
27
28     # Flatten and fully connected layers
29     layers.Flatten(),
30     layers.Dense(128, activation='relu'),
31     layers.Dropout(0.5),
32     layers.Dense(10, activation='softmax')
33 )
34
35 # Model summary
36 model.summary()
37
38 # Compile
39 model.compile(
40     optimizer='adam',
41     loss='sparse_categorical_crossentropy',
42     metrics=['accuracy']
43 )
44
45 # Train
46 history = model.fit(
47     X_train, y_train,
48     epochs=10,
49     batch_size=128,
50     validation_split=0.1,
51     verbose=1
52 )
53
54 # Evaluate
55 test_loss, test_acc = model.evaluate(X_test, y_test, verbose=0)
56 print(f"\nTest accuracy: {test_acc:.4f}")
57
58 # Visualize filters
59 first_layer_weights = model.layers[0].get_weights()[0]
60 fig, axes = plt.subplots(4, 8, figsize=(12, 6))
61 for i, ax in enumerate(axes.flat):
62     if i < 32:
63         ax.imshow(first_layer_weights[:, :, 0, i], cmap='gray')
64         ax.axis('off')
65 plt.suptitle('First Layer Filters')
66 plt.tight_layout()
67 plt.show()

```

5.4 Famous CNN Architectures

5.4.1 LeNet-5 (1998)

One of the earliest CNNs for digit recognition.

5.4.2 AlexNet (2012)

- 5 convolutional layers
- 3 fully connected layers
- ReLU activation

- Dropout regularization
- 60M parameters

5.4.3 VGGNet (2014)

- Very deep (16-19 layers)
- Small 3×3 filters
- Uniform architecture
- 138M parameters

5.4.4 ResNet (2015)

Introduced skip connections to enable very deep networks:

$$\mathbf{H}(\mathbf{x}) = \mathcal{F}(\mathbf{x}) + \mathbf{x} \quad (38)$$

```

1 # ResNet block implementation
2 def residual_block(x, filters, kernel_size=3):
3     # Save input
4     shortcut = x
5
6     # First conv layer
7     x = layers.Conv2D(filters, kernel_size, padding='same')(x)
8     x = layers.BatchNormalization()(x)
9     x = layers.Activation('relu')(x)
10
11    # Second conv layer
12    x = layers.Conv2D(filters, kernel_size, padding='same')(x)
13    x = layers.BatchNormalization()(x)
14
15    # Add shortcut
16    x = layers.Add()([x, shortcut])
17    x = layers.Activation('relu')(x)
18
19    return x
20
21 # Build mini ResNet
22 inputs = layers.Input(shape=(28, 28, 1))
23 x = layers.Conv2D(32, 3, padding='same')(inputs)
24 x = layers.BatchNormalization()(x)
25 x = layers.Activation('relu')(x)
26
27 # Add residual blocks
28 x = residual_block(x, 32)
29 x = residual_block(x, 32)
30
31 x = layers.MaxPooling2D(2)(x)
32 x = layers.Flatten()(x)
33 x = layers.Dense(128, activation='relu')(x)
34 outputs = layers.Dense(10, activation='softmax')(x)
35
36 resnet_model = models.Model(inputs=inputs, outputs=outputs)
37 resnet_model.summary()

```

6 Recurrent Neural Networks (RNNs)

RNNs process sequential data by maintaining hidden states.

6.1 Basic RNN

6.1.1 Mathematical Formulation

At time step t :

$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{b}_h) \quad (39)$$

$$\mathbf{y}_t = \mathbf{W}_{hy}\mathbf{h}_t + \mathbf{b}_y \quad (40)$$

where:

- \mathbf{x}_t : input at time t
- \mathbf{h}_t : hidden state at time t
- \mathbf{y}_t : output at time t
- \mathbf{W}_{hh} , \mathbf{W}_{xh} , \mathbf{W}_{hy} : weight matrices

6.1.2 Backpropagation Through Time (BPTT)

Gradients flow backward through time:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \sum_{t=1}^T \frac{\partial \mathcal{L}_t}{\partial \mathbf{W}} \quad (41)$$

Problem: Vanishing/exploding gradients

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} = \prod_{i=k+1}^t \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \quad (42)$$

6.2 Long Short-Term Memory (LSTM)

LSTMs solve vanishing gradient problem using gating mechanisms.

6.2.1 LSTM Equations

Forget Gate:

$$\mathbf{f}_t = \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f) \quad (43)$$

Input Gate:

$$\mathbf{i}_t = \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \quad (44)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{W}_C[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_C) \quad (45)$$

Cell State Update:

$$\mathbf{C}_t = \mathbf{f}_t \odot \mathbf{C}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{C}}_t \quad (46)$$

Output Gate:

$$\mathbf{o}_t = \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o) \quad (47)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{C}_t) \quad (48)$$

6.3 Gated Recurrent Unit (GRU)

Simplified version of LSTM with fewer parameters.

6.3.1 GRU Equations

Reset Gate:

$$\mathbf{r}_t = \sigma(\mathbf{W}_r[\mathbf{h}_{t-1}, \mathbf{x}_t]) \quad (49)$$

Update Gate:

$$\mathbf{z}_t = \sigma(\mathbf{W}_z[\mathbf{h}_{t-1}, \mathbf{x}_t]) \quad (50)$$

Candidate Hidden State:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}[\mathbf{r}_t \odot \mathbf{h}_{t-1}, \mathbf{x}_t]) \quad (51)$$

Hidden State:

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t \quad (52)$$

6.3.2 Python Implementation

```
1 import tensorflow as tf
2 from tensorflow.keras import layers, models
3 import numpy as np
4
5 # Generate synthetic sequence data
6 def generate_sequence(length=50):
7     X = np.random.randn(1000, length, 1)
8     # Predict sum of sequence
9     y = np.sum(X, axis=1)
10    return X, y
11
12 X_train, y_train = generate_sequence()
13 X_test, y_test = generate_sequence()
14
15 # LSTM Model
16 lstm_model = models.Sequential([
17     layers.LSTM(64, return_sequences=True, input_shape=(50, 1)),
18     layers.Dropout(0.2),
19     layers.LSTM(32),
20     layers.Dropout(0.2),
21     layers.Dense(1)
22 ])
23
24 lstm_model.compile(optimizer='adam', loss='mse', metrics=['mae'])
25 lstm_model.summary()
26
27 # Train
28 history_lstm = lstm_model.fit(
29     X_train, y_train,
30     epochs=20,
31     batch_size=32,
32     validation_split=0.2,
33     verbose=1
34 )
35
36 # GRU Model for comparison
37 gru_model = models.Sequential([
```

```

38     layers.GRU(64, return_sequences=True, input_shape=(50, 1)),
39     layers.Dropout(0.2),
40     layers.GRU(32),
41     layers.Dropout(0.2),
42     layers.Dense(1)
43 )
44
45 gru_model.compile(optimizer='adam', loss='mse', metrics=['mae'])
46
47 history_gru = gru_model.fit(
48     X_train, y_train,
49     epochs=20,
50     batch_size=32,
51     validation_split=0.2,
52     verbose=0
53 )
54
55 # Compare
56 plt.figure(figsize=(10, 5))
57 plt.plot(history_lstm.history['loss'], label='LSTM Train')
58 plt.plot(history_lstm.history['val_loss'], label='LSTM Val')
59 plt.plot(history_gru.history['loss'], label='GRU Train')
60 plt.plot(history_gru.history['val_loss'], label='GRU Val')
61 plt.xlabel('Epoch')
62 plt.ylabel('Loss')
63 plt.title('LSTM vs GRU')
64 plt.legend()
65 plt.grid(True)
66 plt.show()

```

6.4 Bidirectional RNNs

Process sequences in both forward and backward directions:

$$\vec{\mathbf{h}}_t = \text{RNN}(\mathbf{x}_t, \vec{\mathbf{h}}_{t-1}) \quad (53)$$

$$\overleftarrow{\mathbf{h}}_t = \text{RNN}(\mathbf{x}_t, \overleftarrow{\mathbf{h}}_{t+1}) \quad (54)$$

$$\mathbf{h}_t = [\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t] \quad (55)$$

```

1 # Bidirectional LSTM
2 bi_model = models.Sequential([
3     layers.Bidirectional(layers.LSTM(64, return_sequences=True),
4                         input_shape=(50, 1)),
5     layers.Bidirectional(layers.LSTM(32)),
6     layers.Dense(1)
7 ])
8
9 bi_model.compile(optimizer='adam', loss='mse')
10 bi_model.summary()

```

7 Attention Mechanisms

Attention allows models to focus on relevant parts of the input.

7.1 Attention Score

Scaled Dot-Product Attention:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V} \quad (56)$$

where:

- \mathbf{Q} : Query matrix
- \mathbf{K} : Key matrix
- \mathbf{V} : Value matrix
- d_k : dimension of keys

7.2 Multi-Head Attention

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)\mathbf{W}^O \quad (57)$$

where:

$$\text{head}_i = \text{Attention}(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V) \quad (58)$$

```

1 # Simple attention layer implementation
2 class AttentionLayer(layers.Layer):
3     def __init__(self, units):
4         super(AttentionLayer, self).__init__()
5         self.W1 = layers.Dense(units)
6         self.W2 = layers.Dense(units)
7         self.V = layers.Dense(1)
8
9     def call(self, query, values):
10        # query: (batch, hidden)
11        # values: (batch, seq_len, hidden)
12
13        # Expand query to match values shape
14        query_with_time = tf.expand_dims(query, 1)
15
16        # Compute attention scores
17        score = self.V(tf.nn.tanh(
18            self.W1(query_with_time) + self.W2(values)
19        ))
20
21        # Compute attention weights
22        attention_weights = tf.nn.softmax(score, axis=1)
23
24        # Compute context vector
25        context = attention_weights * values
26        context = tf.reduce_sum(context, axis=1)
27
28        return context, attention_weights
29
30 # Example usage in sequence model
31 inputs = layers.Input(shape=(50, 1))
32 lstm_out = layers.LSTM(64, return_sequences=True)(inputs)
33 query = layers.LSTM(64)(inputs)
34 context, weights = AttentionLayer(64)(query, lstm_out)
35 output = layers.Dense(1)(context)

```

```

36
37 attention_model = models.Model(inputs=inputs, outputs=output)
38 attention_model.compile(optimizer='adam', loss='mse')

```

8 Transformers

Transformers rely entirely on attention mechanisms, eliminating recurrence.

8.1 Transformer Architecture

Positional Encoding:

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right) \quad (59)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right) \quad (60)$$

Encoder Block:

$$\mathbf{Z} = \text{MultiHeadAttention}(\mathbf{X}, \mathbf{X}, \mathbf{X}) \quad (61)$$

$$\mathbf{Z}' = \text{LayerNorm}(\mathbf{X} + \mathbf{Z}) \quad (62)$$

$$\mathbf{Y} = \text{FFN}(\mathbf{Z}') \quad (63)$$

$$\mathbf{Y}' = \text{LayerNorm}(\mathbf{Z}' + \mathbf{Y}) \quad (64)$$

```

1 # Simplified Transformer implementation
2 def positional_encoding(length, depth):
3     positions = np.arange(length)[:, np.newaxis]
4     depths = np.arange(depth)[np.newaxis, :] / depth
5
6     angle_rates = 1 / (10000**depths)
7     angle_rads = positions * angle_rates
8
9     pos_encoding = np.zeros(angle_rads.shape)
10    pos_encoding[:, 0::2] = np.sin(angle_rads[:, 0::2])
11    pos_encoding[:, 1::2] = np.cos(angle_rads[:, 1::2])
12
13    return tf.cast(pos_encoding, dtype=tf.float32)
14
15 class TransformerBlock(layers.Layer):
16     def __init__(self, embed_dim, num_heads, ff_dim):
17         super(TransformerBlock, self).__init__()
18         self.att = layers.MultiHeadAttention(
19             num_heads=num_heads, key_dim=embed_dim
20         )
21         self.ffn = keras.Sequential([
22             layers.Dense(ff_dim, activation='relu'),
23             layers.Dense(embed_dim)
24         ])
25         self.layernorm1 = layers.LayerNormalization(epsilon=1e-6)
26         self.layernorm2 = layers.LayerNormalization(epsilon=1e-6)
27         self.dropout1 = layers.Dropout(0.1)
28         self.dropout2 = layers.Dropout(0.1)
29
30     def call(self, inputs, training):
31         attn_output = self.att(inputs, inputs)

```

```

32     attn_output = self.dropout1(attn_output, training=training)
33     out1 = self.layernorm1(inputs + attn_output)
34
35     ffn_output = self.ffn(out1)
36     ffn_output = self.dropout2(ffn_output, training=training)
37     return self.layernorm2(out1 + ffn_output)
38
39 # Build transformer model
40 vocab_size = 10000
41 maxlen = 100
42 embed_dim = 128
43 num_heads = 8
44 ff_dim = 512
45
46 inputs = layers.Input(shape=(maxlen,))
47 embedding_layer = layers.Embedding(vocab_size, embed_dim)(inputs)
48 positions = positional_encoding(maxlen, embed_dim)
49 x = embedding_layer + positions
50
51 transformer_block = TransformerBlock(embed_dim, num_heads, ff_dim)
52 x = transformer_block(x)
53 x = layers.GlobalAveragePooling1D()(x)
54 x = layers.Dropout(0.1)(x)
55 x = layers.Dense(20, activation='relu')(x)
56 outputs = layers.Dense(2, activation='softmax')(x)
57
58 transformer_model = models.Model(inputs=inputs, outputs=outputs)
59 transformer_model.summary()

```

9 Autoencoders

Autoencoders learn compressed representations of data.

9.1 Vanilla Autoencoder

Encoder:

$$\mathbf{z} = f_{\text{enc}}(\mathbf{x}; \boldsymbol{\theta}_{\text{enc}}) \quad (65)$$

Decoder:

$$\hat{\mathbf{x}} = f_{\text{dec}}(\mathbf{z}; \boldsymbol{\theta}_{\text{dec}}) \quad (66)$$

Loss:

$$\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 \quad (67)$$

9.2 Variational Autoencoder (VAE)

VAE learns probabilistic latent representations.

Encoder produces distribution parameters:

$$\boldsymbol{\mu} = f_{\mu}(\mathbf{x}) \quad (68)$$

$$\log \boldsymbol{\sigma}^2 = f_{\sigma}(\mathbf{x}) \quad (69)$$

Sampling (reparameterization trick):

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}) \quad (70)$$

Loss (ELBO):

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - D_{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (71)$$

```

1 # Variational Autoencoder
2 latent_dim = 32
3
4 # Encoder
5 encoder_inputs = layers.Input(shape=(28, 28, 1))
6 x = layers.Conv2D(32, 3, activation='relu', strides=2, padding='same')(encoder_inputs)
7 x = layers.Conv2D(64, 3, activation='relu', strides=2, padding='same')(x)
8 x = layers.Flatten()(x)
9 x = layers.Dense(16, activation='relu')(x)
10
11 z_mean = layers.Dense(latent_dim, name='z_mean')(x)
12 z_log_var = layers.Dense(latent_dim, name='z_log_var')(x)
13
14 # Sampling layer
15 class Sampling(layers.Layer):
16     def call(self, inputs):
17         z_mean, z_log_var = inputs
18         batch = tf.shape(z_mean)[0]
19         dim = tf.shape(z_mean)[1]
20         epsilon = tf.random.normal(shape=(batch, dim))
21         return z_mean + tf.exp(0.5 * z_log_var) * epsilon
22
23 z = Sampling()([z_mean, z_log_var])
24
25 encoder = models.Model(encoder_inputs, [z_mean, z_log_var, z], name='encoder')
26
27 # Decoder
28 latent_inputs = layers.Input(shape=(latent_dim,))
29 x = layers.Dense(7 * 7 * 64, activation='relu')(latent_inputs)
30 x = layers.Reshape((7, 7, 64))(x)
31 x = layers.Conv2DTranspose(64, 3, activation='relu', strides=2, padding='same')(x)
32 x = layers.Conv2DTranspose(32, 3, activation='relu', strides=2, padding='same')(x)
33 decoder_outputs = layers.Conv2DTranspose(1, 3, activation='sigmoid',
34                                         padding='same')(x)
35
36 decoder = models.Model(latent_inputs, decoder_outputs, name='decoder')
37
38 # VAE model
39 class VAE(keras.Model):
40     def __init__(self, encoder, decoder, **kwargs):
41         super(VAE, self).__init__(**kwargs)
42         self.encoder = encoder
43         self.decoder = decoder
44         self.total_loss_tracker = keras.metrics.Mean(name='total_loss')
45         self.reconstruction_loss_tracker = keras.metrics.Mean(name='reconstruction_loss')
46         self.kl_loss_tracker = keras.metrics.Mean(name='kl_loss')
47
48     @property
49     def metrics(self):

```

```

49     return [
50         self.total_loss_tracker,
51         self.reconstruction_loss_tracker,
52         self.kl_loss_tracker
53     ]
54
55     def train_step(self, data):
56         with tf.GradientTape() as tape:
57             z_mean, z_log_var, z = self.encoder(data)
58             reconstruction = self.decoder(z)
59
60             reconstruction_loss = tf.reduce_mean(
61                 tf.reduce_sum(
62                     keras.losses.binary_crossentropy(data,
63                         reconstruction),
64                     axis=(1, 2)
65                 )
66             )
67
68             kl_loss = -0.5 * (1 + z_log_var - tf.square(z_mean) - tf.
69                 exp(z_log_var))
70             kl_loss = tf.reduce_mean(tf.reduce_sum(kl_loss, axis=1))
71
72             total_loss = reconstruction_loss + kl_loss
73
74             grads = tape.gradient(total_loss, self.trainable_weights)
75             self.optimizer.apply_gradients(zip(grads, self.
76                 trainable_weights))
77
78             self.total_loss_tracker.update_state(total_loss)
79             self.reconstruction_loss_tracker.update_state(
80                 reconstruction_loss)
81             self.kl_loss_tracker.update_state(kl_loss)
82
83         return {
84             'loss': self.total_loss_tracker.result(),
85             'reconstruction_loss': self.reconstruction_loss_tracker.
86                 result(),
87             'kl_loss': self.kl_loss_tracker.result()
88         }
89
90     vae = VAE(encoder, decoder)
91     vae.compile(optimizer='adam')
92
93 # Train on MNIST
94 (x_train, _), (x_test, _) = keras.datasets.mnist.load_data()
95 x_train = np.expand_dims(x_train, -1).astype('float32') / 255.0
96 x_test = np.expand_dims(x_test, -1).astype('float32') / 255.0
97
98 vae.fit(x_train, epochs=30, batch_size=128)
99
100 # Generate new images
101 n = 15
102 digit_size = 28
103 figure = np.zeros((digit_size * n, digit_size * n))
104 grid_x = np.linspace(-3, 3, n)
105 grid_y = np.linspace(-3, 3, n)[::-1]

```

```

102
103 for i, yi in enumerate(grid_y):
104     for j, xi in enumerate(grid_x):
105         z_sample = np.array([[xi, yi] + [0] * (latent_dim - 2)])
106         x_decoded = decoder.predict(z_sample, verbose=0)
107         digit = x_decoded[0].reshape(digit_size, digit_size)
108         figure[i * digit_size: (i + 1) * digit_size,
109                j * digit_size: (j + 1) * digit_size] = digit
110
111 plt.figure(figsize=(10, 10))
112 plt.imshow(figure, cmap='Greys_r')
113 plt.title('VAE Generated Digits')
114 plt.axis('off')
115 plt.show()

```

10 Generative Adversarial Networks (GANs)

GANs consist of two networks competing in a minimax game.

10.1 GAN Objective

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log(1 - D(G(\mathbf{z})))] \quad (72)$$

where:

- G : Generator network
- D : Discriminator network
- \mathbf{z} : Random noise vector

```

1 # Simple GAN implementation
2 latent_dim = 100
3
4 # Generator
5 generator = models.Sequential([
6     layers.Dense(7 * 7 * 128, input_dim=latent_dim),
7     layers.Reshape((7, 7, 128)),
8     layers.BatchNormalization(),
9     layers.LeakyReLU(alpha=0.2),
10
11    layers.Conv2DTranspose(128, 4, strides=2, padding='same'),
12    layers.BatchNormalization(),
13    layers.LeakyReLU(alpha=0.2),
14
15    layers.Conv2DTranspose(64, 4, strides=2, padding='same'),
16    layers.BatchNormalization(),
17    layers.LeakyReLU(alpha=0.2),
18
19    layers.Conv2D(1, 3, padding='same', activation='tanh')
20], name='generator')
21
22 # Discriminator
23 discriminator = models.Sequential([
24     layers.Conv2D(64, 3, strides=2, padding='same',
25                  input_shape=(28, 28, 1)),
26     layers.LeakyReLU(alpha=0.2),

```

```

27     layers.Dropout(0.3),
28
29     layers.Conv2D(128, 3, strides=2, padding='same'),
30     layers.LeakyReLU(alpha=0.2),
31     layers.Dropout(0.3),
32
33     layers.Flatten(),
34     layers.Dense(1, activation='sigmoid')
35 ], name='discriminator')
36
37 # GAN model
38 class GAN(keras.Model):
39     def __init__(self, discriminator, generator, latent_dim):
40         super(GAN, self).__init__()
41         self.discriminator = discriminator
42         self.generator = generator
43         self.latent_dim = latent_dim
44
45     def compile(self, d_optimizer, g_optimizer, loss_fn):
46         super(GAN, self).compile()
47         self.d_optimizer = d_optimizer
48         self.g_optimizer = g_optimizer
49         self.loss_fn = loss_fn
50
51     def train_step(self, real_images):
52         batch_size = tf.shape(real_images)[0]
53
54         # Generate fake images
55         random_latent_vectors = tf.random.normal(
56             shape=(batch_size, self.latent_dim)
57         )
58         generated_images = self.generator(random_latent_vectors)
59
60         # Combine with real images
61         combined_images = tf.concat([generated_images, real_images],
62                                     axis=0)
63         labels = tf.concat([
64             tf.ones((batch_size, 1)),
65             tf.zeros((batch_size, 1))
66         ], axis=0)
67
68         # Add noise to labels
69         labels += 0.05 * tf.random.uniform(tf.shape(labels))
70
71         # Train discriminator
72         with tf.GradientTape() as tape:
73             predictions = self.discriminator(combined_images)
74             d_loss = self.loss_fn(labels, predictions)
75
76             grads = tape.gradient(d_loss, self.discriminator.
77                                   trainable_weights)
78             self.d_optimizer.apply_gradients(
79                 zip(grads, self.discriminator.trainable_weights)
80             )
81
82         # Train generator
83         random_latent_vectors = tf.random.normal(
84             shape=(batch_size, self.latent_dim)

```

```

83     )
84     misleading_labels = tf.zeros((batch_size, 1))
85
86     with tf.GradientTape() as tape:
87         predictions = self.discriminator(
88             self.generator(random_latent_vectors)
89         )
90         g_loss = self.loss_fn(misleading_labels, predictions)
91
92     grads = tape.gradient(g_loss, self.generator.trainable_weights)
93     self.g_optimizer.apply_gradients(
94         zip(grads, self.generator.trainable_weights)
95     )
96
97     return {'d_loss': d_loss, 'g_loss': g_loss}
98
99 # Create and compile GAN
100 gan = GAN(discriminator=discriminator,
101             generator=generator,
102             latent_dim=latent_dim)
103
104 gan.compile(
105     d_optimizer=keras.optimizers.Adam(learning_rate=0.0002, beta_1=0.5),
106     ,
107     g_optimizer=keras.optimizers.Adam(learning_rate=0.0002, beta_1=0.5),
108     ,
109     loss_fn=keras.losses.BinaryCrossentropy()
110 )
111
112 # Train
113 (x_train, _), _ = keras.datasets.mnist.load_data()
114 x_train = x_train.reshape(-1, 28, 28, 1).astype('float32')
115 x_train = (x_train - 127.5) / 127.5 # Normalize to [-1, 1]
116
117 gan.fit(x_train, epochs=50, batch_size=128)
118
119 # Generate images
120 noise = tf.random.normal([16, latent_dim])
121 generated_images = generator(noise, training=False)
122
123 fig, axes = plt.subplots(4, 4, figsize=(8, 8))
124 for i, ax in enumerate(axes.flat):
125     ax.imshow(generated_images[i, :, :, 0], cmap='gray')
126     ax.axis('off')
127 plt.suptitle('GAN Generated Images')
128 plt.tight_layout()
129 plt.show()

```

11 Transfer Learning

Transfer learning leverages pre-trained models for new tasks.

11.1 Approaches

1. **Feature Extraction:** Freeze pre-trained layers, train new classifier
2. **Fine-Tuning:** Unfreeze some layers and retrain

3. **Full Training:** Use pre-trained weights as initialization

```
1 # Transfer learning with VGG16
2 from tensorflow.keras.applications import VGG16
3 from tensorflow.keras.preprocessing.image import ImageDataGenerator
4
5 # Load pre-trained VGG16 (without top layers)
6 base_model = VGG16(
7     weights='imagenet',
8     include_top=False,
9     input_shape=(224, 224, 3)
10)
11
12 # Freeze base model
13 base_model.trainable = False
14
15 # Add custom classifier
16 inputs = keras.Input(shape=(224, 224, 3))
17 x = base_model(inputs, training=False)
18 x = layers.GlobalAveragePooling2D()(x)
19 x = layers.Dense(256, activation='relu')(x)
20 x = layers.Dropout(0.5)(x)
21 outputs = layers.Dense(10, activation='softmax')(x)
22
23 model = keras.Model(inputs, outputs)
24
25 model.compile(
26     optimizer='adam',
27     loss='sparse_categorical_crossentropy',
28     metrics=['accuracy']
29)
30
31 model.summary()
32
33 # Fine-tuning: unfreeze last few layers
34 base_model.trainable = True
35 for layer in base_model.layers[:-4]:
36     layer.trainable = False
37
38 model.compile(
39     optimizer=keras.optimizers.Adam(1e-5), # Lower learning rate
40     loss='sparse_categorical_crossentropy',
41     metrics=['accuracy']
42)
```

12 Best Practices & Tips

12.1 Training Strategies

1. Learning Rate Scheduling

$$\alpha_t = \alpha_0 \cdot \text{decay_factor}^{t/\text{decay_steps}} \quad (73)$$

2. **Early Stopping:** Monitor validation loss, stop if no improvement

3. **Gradient Clipping:** Prevent exploding gradients

$$\mathbf{g} := \min\left(1, \frac{\text{threshold}}{\|\mathbf{g}\|}\right) \cdot \mathbf{g} \quad (74)$$

4. **Data Augmentation:** Artificially expand training data

5. **Batch Size Selection:** Larger batches = faster but less generalization

```

1 # Learning rate scheduling
2 from tensorflow.keras.callbacks import ReduceLROnPlateau, EarlyStopping
3
4 reduce_lr = ReduceLROnPlateau(
5     monitor='val_loss',
6     factor=0.5,
7     patience=5,
8     min_lr=1e-7,
9     verbose=1
10)
11
12 early_stop = EarlyStopping(
13     monitor='val_loss',
14     patience=10,
15     restore_best_weights=True,
16     verbose=1
17)
18
19 # Data augmentation for images
20 datagen = ImageDataGenerator(
21     rotation_range=15,
22     width_shift_range=0.1,
23     height_shift_range=0.1,
24     horizontal_flip=True,
25     zoom_range=0.1,
26     shear_range=0.1,
27     fill_mode='nearest'
28)
29
30 # Training with callbacks
31 history = model.fit(
32     datagen.flow(X_train, y_train, batch_size=32),
33     epochs=100,
34     validation_data=(X_val, y_val),
35     callbacks=[reduce_lr, early_stop]
36)

```

12.2 Model Evaluation

Classification Metrics:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (75)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (76)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (77)$$

$$\text{F1-Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (78)$$

```

1  from sklearn.metrics import classification_report, confusion_matrix
2  import seaborn as sns
3
4  # Predictions
5  y_pred = model.predict(X_test)
6  y_pred_classes = np.argmax(y_pred, axis=1)
7
8  # Classification report
9  print(classification_report(y_test, y_pred_classes))
10
11 # Confusion matrix
12 cm = confusion_matrix(y_test, y_pred_classes)
13 plt.figure(figsize=(10, 8))
14 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
15 plt.xlabel('Predicted')
16 plt.ylabel('Actual')
17 plt.title('Confusion Matrix')
18 plt.show()

```

12.3 Hyperparameter Tuning

```

1  from keras_tuner import RandomSearch
2
3  def build_model(hp):
4      model = keras.Sequential()
5
6      # Tune number of layers
7      for i in range(hp.Int('num_layers', 1, 4)):
8          model.add(layers.Dense(
9              units=hp.Int(f'units_{i}', 32, 512, step=32),
10             activation='relu'))
11
12          model.add(layers.Dropout(
13              hp.Float(f'dropout_{i}', 0, 0.5, step=0.1)))
14
15
16      model.add(layers.Dense(10, activation='softmax'))
17
18      # Tune learning rate
19      model.compile(
20          optimizer=keras.optimizers.Adam(
21              hp.Float('learning_rate', 1e-4, 1e-2, sampling='log')
22          ),
23          loss='sparse_categorical_crossentropy',
24          metrics=['accuracy']
25      )
26
27
28      return model
29
30
31 tuner = RandomSearch(
32     build_model,
33     objective='val_accuracy',
34     max_trials=10,
35     directory='tuning',
36     project_name='mnist'
37 )

```

```

36
37 tuner.search(X_train, y_train, epochs=10, validation_split=0.2)
38 best_model = tuner.get_best_models(num_models=1)[0]

```

12.4 Debugging Neural Networks

Common Issues:

- **Vanishing Gradients:** Use ReLU, batch normalization, skip connections
- **Exploding Gradients:** Use gradient clipping, lower learning rate
- **Overfitting:** Add dropout, L2 regularization, more data
- **Underfitting:** Increase model capacity, train longer
- **Dead ReLU:** Use Leaky ReLU, lower learning rate

```

1 # Gradient monitoring
2 class GradientLogger(keras.callbacks.Callback):
3     def on_epoch_end(self, epoch, logs=None):
4         for layer in self.model.layers:
5             if hasattr(layer, 'kernel'):
6                 weights = layer.get_weights()[0]
7                 print(f'{layer.name}: mean={np.mean(weights):.4f}, '
8                      f'std={np.std(weights):.4f}')
9
10 # Check for dead neurons
11 def check_dead_neurons(model, X_sample):
12     activations = []
13     temp_model = keras.Model(
14         inputs=model.input,
15         outputs=[layer.output for layer in model.layers]
16     )
17
18     outputs = temp_model.predict(X_sample)
19
20     for i, output in enumerate(outputs):
21         if len(output.shape) == 2: # Dense layer
22             dead = np.sum(output == 0, axis=0) / output.shape[0]
23             print(f'Layer {i}: {np.sum(dead > 0.9)} dead neurons '
24                  f'({100*np.mean(dead > 0.9):.1f}%)')
25
26 check_dead_neurons(model, X_train[:1000])

```

13 Advanced Topics

13.1 Neural Architecture Search (NAS)

Automatically discover optimal architectures using:

- Reinforcement Learning
- Evolutionary Algorithms
- Gradient-based methods (DARTS)

13.2 Meta-Learning

Learning to learn: models that adapt quickly to new tasks with few examples.

MAML (Model-Agnostic Meta-Learning):

$$\boldsymbol{\theta}^* = \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{T}_i}(f_{\boldsymbol{\theta}}) \quad (79)$$

13.3 Neural ODEs

Model dynamics as continuous transformations:

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \boldsymbol{\theta}) \quad (80)$$

13.4 Graph Neural Networks (GNNs)

Process graph-structured data:

$$\mathbf{h}_v^{(k+1)} = \sigma \left(\mathbf{W}^{(k)} \mathbf{h}_v^{(k)} + \sum_{u \in \mathcal{N}(v)} \mathbf{h}_u^{(k)} \right) \quad (81)$$

14 Practical Projects

14.1 Image Classification

```
1 # Complete image classification pipeline
2 import tensorflow as tf
3 from tensorflow.keras.applications import EfficientNetB0
4
5 # Load and preprocess data
6 def load_dataset(data_dir):
7     train_ds = tf.keras.preprocessing.image_dataset_from_directory(
8         data_dir + '/train',
9         validation_split=0.2,
10        subset='training',
11        seed=123,
12        image_size=(224, 224),
13        batch_size=32
14    )
15
16    val_ds = tf.keras.preprocessing.image_dataset_from_directory(
17        data_dir + '/train',
18        validation_split=0.2,
19        subset='validation',
20        seed=123,
21        image_size=(224, 224),
22        batch_size=32
23    )
24
25    return train_ds, val_ds
26
27 # Build model with transfer learning
28 def build_classifier(num_classes):
29     base_model = EfficientNetB0(
30         include_top=False,
31         weights='imagenet',
```

```

32         input_shape=(224, 224, 3)
33     )
34     base_model.trainable = False
35
36     inputs = keras.Input(shape=(224, 224, 3))
37     x = tf.keras.applications.efficientnet.preprocess_input(inputs)
38     x = base_model(x, training=False)
39     x = layers.GlobalAveragePooling2D()(x)
40     x = layers.Dropout(0.5)(x)
41     outputs = layers.Dense(num_classes, activation='softmax')(x)
42
43     model = keras.Model(inputs, outputs)
44     return model
45
46 # Train with mixed precision
47 from tensorflow.keras import mixed_precision
48 mixed_precision.set_global_policy('mixed_float16')
49
50 model = build_classifier(num_classes=10)
51 model.compile(
52     optimizer='adam',
53     loss='sparse_categorical_crossentropy',
54     metrics=['accuracy']
55 )
56
57 # Callbacks
58 callbacks = [
59     keras.callbacks.ModelCheckpoint('best_model.h5', save_best_only=
60         True),
61     keras.callbacks.ReduceLROnPlateau(patience=3),
62     keras.callbacks.EarlyStopping(patience=10),
63     keras.callbacks.TensorBoard(log_dir='./logs')
64 ]
65
66 # Train
67 history = model.fit(
68     train_ds,
69     validation_data=val_ds,
70     epochs=50,
71     callbacks=callbacks
72 )

```

14.2 Text Generation with RNN

```

1 # Character-level text generation
2 class TextGenerator:
3     def __init__(self, text, seq_length=100):
4         self.text = text
5         self.seq_length = seq_length
6         self.chars = sorted(list(set(text)))
7         self.char_to_idx = {c: i for i, c in enumerate(self.chars)}
8         self.idx_to_char = {i: c for i, c in enumerate(self.chars)}
9
10    def prepare_data(self):
11        sequences = []
12        next_chars = []
13

```

```

14     for i in range(0, len(self.text) - self.seq_length):
15         sequences.append(self.text[i:i + self.seq_length])
16         next_chars.append(self.text[i + self.seq_length])
17
18     X = np.zeros((len(sequences), self.seq_length, len(self.chars)))
19     y = np.zeros((len(sequences), len(self.chars)))
20
21     for i, seq in enumerate(sequences):
22         for t, char in enumerate(seq):
23             X[i, t, self.char_to_idx[char]] = 1
24             y[i, self.char_to_idx[next_chars[i]]] = 1
25
26     return X, y
27
28 def build_model(self):
29     model = keras.Sequential([
30         layers.LSTM(128, input_shape=(self.seq_length, len(self.
31             chars)),
32                     return_sequences=True),
33         layers.Dropout(0.2),
34         layers.LSTM(128),
35         layers.Dropout(0.2),
36         layers.Dense(len(self.chars), activation='softmax')
37     ])
38
39     model.compile(
40         optimizer='adam',
41         loss='categorical_crossentropy',
42         metrics=['accuracy']
43     )
44
45     return model
46
47 def generate_text(self, model, seed_text, length=400, temperature
48 =1.0):
49     generated = seed_text
50
51     for _ in range(length):
52         x = np.zeros((1, self.seq_length, len(self.chars)))
53         for t, char in enumerate(seed_text):
54             x[0, t, self.char_to_idx[char]] = 1
55
56         preds = model.predict(x, verbose=0)[0]
57         preds = np.log(preds + 1e-7) / temperature
58         exp_preds = np.exp(preds)
59         preds = exp_preds / np.sum(exp_preds)
60
61         next_idx = np.random.choice(len(self.chars), p=preds)
62         next_char = self.idx_to_char[next_idx]
63
64         generated += next_char
65         seed_text = seed_text[1:] + next_char
66
67     return generated
68
69 # Usage
70 with open('shakespeare.txt', 'r') as f:

```

```

69     text = f.read().lower()
70
71 gen = TextGenerator(text, seq_length=100)
72 X, y = gen.prepare_data()
73 model = gen.build_model()
74 model.fit(X, y, batch_size=128, epochs=50)
75
76 # Generate text
77 seed = text[1000:1100]
78 generated = gen.generate_text(model, seed, length=500, temperature=0.5)
79 print(generated)

```

15 Resources & Further Reading

15.1 Books

- Deep Learning by Goodfellow, Bengio, and Courville
- Neural Networks and Deep Learning by Michael Nielsen
- Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

15.2 Online Courses

- Stanford CS231n: Convolutional Neural Networks
- Stanford CS224n: Natural Language Processing
- Fast.ai Practical Deep Learning
- DeepLearning.AI Specialization

15.3 Research Papers

- ImageNet Classification with Deep CNNs (AlexNet)
- Deep Residual Learning (ResNet)
- Attention Is All You Need (Transformer)
- Generative Adversarial Networks (GANs)
- Auto-Encoding Variational Bayes (VAE)

16 Conclusion

Deep learning has revolutionized artificial intelligence, enabling breakthroughs in computer vision, natural language processing, and beyond. Key takeaways:

1. **Architecture matters:** Choose appropriate architectures (CNNs for images, RNNs/-Transformers for sequences)
2. **Regularization is crucial:** Use dropout, batch normalization, and data augmentation
3. **Optimization techniques:** Adam is often a good default, but experiment with others
4. **Transfer learning:** Leverage pre-trained models when possible

5. Hyperparameter tuning: Systematic search can significantly improve performance

6. Monitor training: Use visualization and callbacks to detect issues early

The field continues to evolve rapidly with new architectures, training techniques, and applications emerging constantly. Stay curious, experiment, and keep learning!