# Neural Network Framework Implementation

A Step-by-Step Guide to Building Your Own Deep Learning Library

Morteza Maleki

Tarbiat Modares University

May 21, 2025

# Table of Contents

# Overview of the NPDL Framework

- A NumPy-based Deep Learning Framework
- Modular design with separate components:
  - Activation functions
  - Weight initializers
  - Neural network layers
  - Loss functions
  - Model construction
  - Optimizers
- Implementation follows object-oriented principles
- Designed for educational purposes

# Backpropagation - The Core of Neural Networks

- Backpropagation: The algorithm that powers neural network learning
- Key components of a neural network training loop:
  1. **Forward pass**: Compute predictions
  2. **Loss calculation**: Measure error
  3. **Backward pass**: Compute gradients
  4. **Parameter update**: Improve model
- Focus of this section: Understanding how gradients flow backward through the network
- Goal: Connect mathematical theory to our code implementation
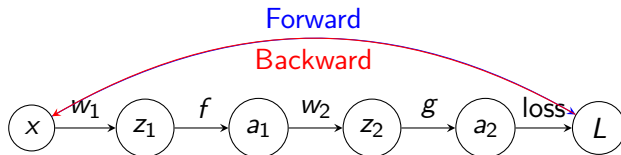
# The Chain Rule - Foundation of Backpropagation

- Neural networks are composed of nested functions
- Chain rule from calculus: If $y = f(g(x))$, then $\frac{dy}{dx} = \frac{dy}{dg} \cdot \frac{dg}{dx}$
- For neural networks with many layers:

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial y_k} \cdot \frac{\partial y_k}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_{ij}}$$

Where:
- $L$ is the loss
- $y_k$ is the output
- $a_j$ is the activation
- $z_j$ is the pre-activation
- $w_{ij}$ is a weight parameter

# Computational Graph Perspective



- Forward pass (blue): Compute outputs from inputs
- Backward pass (red): Propagate gradients from outputs to inputs
- Each node knows:
  - How to compute its output (forward)
  - How to compute gradients w.r.t. its inputs (backward)

## A 2-Layer Neural Network Example

Consider a simple 2-layer neural network:

$$z^{[1]} = W^{[1]}x + b^{[1]} \tag{1}$$

$$a^{[1]} = f(z^{[1]}) \tag{2}$$

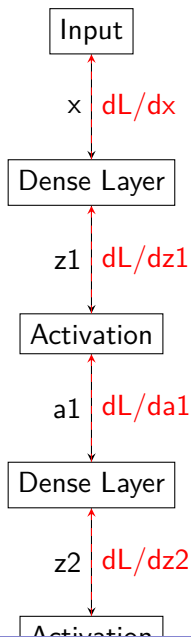$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \tag{3}$$

$$\hat{y} = g(z^{[2]}) \tag{4}$$

$$L = \text{loss}(\hat{y}, y) \tag{5}$$

- Computing gradients using the chain rule:

$$\frac{\partial L}{\partial W^{[2]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial W^{[2]}} \tag{6}$$

$$\frac{\partial L}{\partial W^{[1]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial W^{[1]}} \tag{7}$$

# Gradient Flow in Our Framework

Input

x  |  dL/dx

Dense Layer

z1  |  dL/dz1

Activation

a1  |  dL/da1

Dense Layer

z2  |  dL/dz2

Activation

# The Layer Interface in Our Framework

- Each layer in our framework follows the same interface:
  1. `forward()` - Computes layer output
  2. `backward()` - Computes gradients
  3. Stores parameters in `self.params`
  4. Stores gradients in `self.grads`
- This consistent interface enables automatic differentiation
- Models chain these layers together in both passes:
  - Forward: from first layer to last
  - Backward: from last layer to first

# Example: Gradient Flow in a Dense Layer

```python
1  class Dense(Layer):
2      def forward(self, inputs):
3          """Forward pass: y = Wx + b"""
4          self.inputs = inputs  # Store for backward pass
5          return np.dot(inputs, self.weights) + self.bias
6
7      def backward(self, output_grad):
8          """Backward pass: Compute gradients"""
9          # Compute weight gradients (L/W)
10         self.grads['weights'] = np.dot(self.inputs.T, output_grad)
11
12         # Compute bias gradients (L/b)
13         self.grads['bias'] = np.sum(output_grad, axis=0, keepdims=True)
14
15         # Compute input gradients (L/x) to pass to previous layer
16         return np.dot(output_grad, self.weights.T)
17
```

- Dense layer implements matrix multiplication: $y = Wx + b$
- Backward pass uses calculus rules:
  - $\frac{\partial L}{\partial W} = x^T \cdot \frac{\partial L}{\partial y}$
  - $\frac{\partial L}{\partial b} = \sum \frac{\partial L}{\partial y}$
  - $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot W^T$

# Example: Gradient Flow in an Activation Function

```python
1  class ReLU(Activation):
2      def forward(self, x):
3          """Forward pass: f(x) = max(0, x)"""
4          self.input = x   # Store for backward pass
5          self.output = np.maximum(0, x)
6          return self.output
7
8      def backward(self, output_grad):
9          """Backward pass: Compute gradients"""
10         # Derivative of ReLU: 1 if x > 0, else 0
11         relu_grad = (self.input > 0).astype(float)
12
13         # Chain rule: multiply upstream gradient with local gradient
14         return output_grad * relu_grad
15
```

- Activation functions apply element-wise operations
- Backward pass computes local derivatives and applies chain rule
- For ReLU: $\frac{d}{dx} \max(0, x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$

# How Backward Passes Connect in the Framework

- Sequential model chains the backward passes automatically
- The flow of gradients:
  1. Loss function computes initial gradient: $\frac{\partial L}{\partial \hat{y}}$
  2. Each layer receives gradient w.r.t. its output: $\frac{\partial L}{\partial \text{out}}$
  3. Each layer computes:
     - Gradients for its parameters: $\frac{\partial L}{\partial \theta}$
     - Gradient w.r.t. its input: $\frac{\partial L}{\partial \text{in}}$
  4. Gradient w.r.t. input gets passed to previous layer
- Key insight: Each layer only needs to compute local gradients
- The framework handles chaining the gradients together

# Sequential Model's Backward Pass

```
1  class Sequential(object):
2      def __init__(self, layers=None):
3          self.layers = layers if layers is not None else []
4
5      def forward(self, inputs):
6          """Forward pass through all layers in sequence"""
7          for layer in self.layers:
8              inputs = layer.forward(inputs)
9          return inputs
10
11     def backward(self, grad):
12         """Backward pass through all layers in reverse"""
13         for layer in reversed(self.layers):
14             grad = layer.backward(grad)
15         return grad
16
```

- `Sequential` model chains layers together
- **Forward**: Process inputs through layers in order
- **Backward**: Process gradients through layers in reverse order
- Each layer receives the gradient from the layer ahead
- The full chain rule is automatically implemented

# Implementation to Mathematics Mapping

| Component | Forward | Backward |
|-----------|---------|----------|
| Dense Layer | $z = Wx + b$ | $\frac{\partial L}{\partial W} = x^T \frac{\partial L}{\partial z}$ <br> $\frac{\partial L}{\partial b} = \sum \frac{\partial L}{\partial z}$ <br> $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} W^T$ |
| ReLU | $a = \max(0, z)$ | $\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \mathbf{1}_{z>0}$ |
| Sigmoid | $a = \frac{1}{1+e^{-z}}$ | $\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot a(1-a)$ |
| Softmax | $a_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$ | Complex Jacobian matrix |
| MSE Loss | $L = \frac{1}{2}(y - \hat{y})^2$ | $\frac{\partial L}{\partial \hat{y}} = \hat{y} - y$ |

- Each component has a mathematical operation and its derivative
- The code implementation directly follows these equations
- Understanding the math makes the code implementation clear

# Updating Parameters with Optimizers

- After computing gradients, optimizers update parameters
- Stochastic Gradient Descent (SGD):

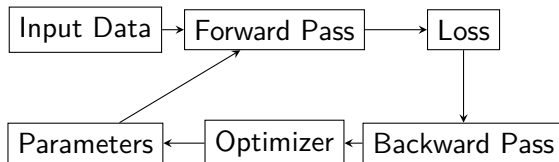$$\theta_{t+1} = \theta_t - \alpha \nabla_\theta J(\theta)$$

- SGD with momentum:

$$v_{t+1} = \gamma v_t + \alpha \nabla_\theta J(\theta) \tag{8}$$
$$\theta_{t+1} = \theta_t - v_{t+1} \tag{9}$$

- Each optimizer implements the `update(layer)` method:
  - Reads gradients from `layer.grads`
  - Updates parameters in `layer.params`
  - May maintain its own state (e.g., momentum)

# The Complete Training Loop

1. **Forward pass**: Input $\rightarrow$ Model $\rightarrow$ Prediction
2. **Loss calculation**: Compare prediction with target
3. **Backward pass**:
   - Start with loss gradient
   - Propagate through model in reverse
   - Collect parameter gradients
4. **Parameter update**:
   - Apply optimizer to update weights using gradients
5. Repeat for many examples and epochs

# Summary: Backpropagation in Our Framework

- Key components implementing backpropagation:
  - **Layers**: Compute outputs and gradients
  - **Activations**: Add non-linearity and corresponding gradients
  - **Loss functions**: Measure error and provide initial gradient
  - **Optimizers**: Update parameters using gradients
  - **Models**: Chain components together
- Design principles:
  - Modularity: Each component has specific responsibility
  - Unified interface: `forward()`/`backward()` methods
  - Mathematical correspondence: Code follows math directly
  - Automatic differentiation: Chain rule applied automatically
- Next, we'll explore each component in detail

# Activation Functions - Overview

- Activation functions introduce non-linearity
- Base abstract class with common interface
- Implementations of common activation functions:
  - Sigmoid
  - Tanh
  - ReLU
  - Softmax
- Each has forward and backward methods

# Activation Functions - Base Class

```python
class Activation(object):
    """Base class for all activation functions."""

    def forward(self, x):
        """Forward pass."""
        raise NotImplementedError

    def backward(self, output_grad):
        """Backward pass."""
        raise NotImplementedError
```

# Activation Functions - Sigmoid

```
1  class Sigmoid(Activation):
2      """Sigmoid activation function."""
3
4      def forward(self, x):
5          """Forward pass for sigmoid function.
6
7          Args:
8              x: Input numpy array.
9          """
10         self.output = 1.0 / (1.0 + np.exp(-x))
11         return self.output
12
13     def backward(self, output_grad):
14         """Backward pass for sigmoid function.
15
16         Args:
17             output_grad: Gradient of the cost with respect to the
       output.
18         """
19         return output_grad * self.output * (1 - self.output)
```

# Sigmoid Activation Function - Mathematical Definition

- The sigmoid function is defined as:

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

- Key properties:
  - Outputs values between 0 and 1
  - Smooth and differentiable everywhere
  - Historically popular but prone to vanishing gradient
  - S-shaped curve (sigmoid shape)
- Used primarily in:
  - Binary classification (output layer)
  - Early neural networks (mostly replaced by ReLU)
  - Gates in recurrent neural networks (LSTM, GRU)

# Sigmoid Activation Function - Gradient Derivation

$$\sigma(x) = \frac{1}{1 + e^{-x}} \tag{10}$$

$$\frac{d\sigma(x)}{dx} = \frac{d}{dx}\left(\frac{1}{1 + e^{-x}}\right) \tag{11}$$

$$= \frac{d}{dx}(1 + e^{-x})^{-1} \tag{12}$$

$$= -(1 + e^{-x})^{-2} \cdot \frac{d}{dx}(1 + e^{-x}) \tag{13}$$

$$= -(1 + e^{-x})^{-2} \cdot (-e^{-x}) \tag{14}$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2} \tag{15}$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \tag{16}$$

$$= \frac{1}{1 + e^{-x}} \cdot \left(1 - \frac{1}{1 + e^{-x}}\right) \tag{17}$$

# Sigmoid Activation Function - Implementation Details

- Forward pass:
  - Calculate $\sigma(x) = \frac{1}{1+e^{-x}}$
  - Store output for use in backward pass
  - Handle numerical stability with clipping for extreme values
- Backward pass:
  - Use the elegant form of gradient: $\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$
  - Multiply input gradient by this derivative (chain rule)
  - No need to reference original input - only output is needed
- Implementation challenges:
  - Saturation for large positive/negative inputs
  - Vanishing gradient problem when chaining multiple sigmoids

# ReLU Activation Function - Detailed Overview

- Rectified Linear Unit defined as:

$$\text{ReLU}(x) = \max(0, x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

- Key properties:
  - Simple, computationally efficient
  - Non-linear despite simple form
  - Sparse activation (many neurons output zero)
  - No vanishing gradient for positive inputs
  - Allows for deeper networks
- Limitations:
  - "Dying ReLU" problem when neurons get stuck at 0
  - Non-zero centered outputs
  - Unbounded positive activation

# ReLU Activation Function - Gradient Analysis

$$\frac{d\text{ReLU}(x)}{dx} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \\ \text{undefined} & \text{if } x = 0 \end{cases} \tag{19}$$

- Gradient is either 0 or 1 (easy to compute)
- No saturation for positive values (solves vanishing gradient)
- Promotes sparsity in the network
- Gradient at $x = 0$ is technically undefined, but usually set to 0 or 1
- Computationally efficient gradient - just check if input was positive

# ReLU Variants - LeakyReLU, PReLU, ELU

- **LeakyReLU**: Allows small negative values with a fixed slope

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0 \end{cases}$$

  where $\alpha$ is a small constant (e.g., 0.01)

- **Parametric ReLU (PReLU)**: Learns the slope parameter $\alpha$ during training

- **Exponential Linear Unit (ELU)**:

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

- These variants help address the "dying ReLU" problem while maintaining ReLU's advantages

# Activation Functions - ReLU

```python
1  class ReLU(Activation):
2      """ReLU activation function."""
3
4      def forward(self, x):
5          """Forward pass for ReLU function.
6
7          Args:
8              x: Input numpy array.
9          """
10         self.input = x
11         self.output = np.maximum(x, 0)
12         return self.output
13
14     def backward(self, output_grad):
15         """Backward pass for ReLU function.
16
17         Args:
18             output_grad: Gradient of the cost with respect to the
       output.
19         """
20         return output_grad * (self.input > 0)
```

# Softmax Activation Function - Comprehensive Overview

- Softmax converts a vector of values into a probability distribution:

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^{n} e^{x_j}}$$

- Key properties:
  - Outputs are in range [0,1] and sum to 1
  - Preserves relative order of inputs (monotonic)
  - Emphasizes largest values, suppresses smaller ones
  - Not element-wise (depends on all input values)
- Use cases:
  - Output layer for multi-class classification
  - Attention mechanisms in transformers
  - Any scenario requiring normalized probabilities

# Softmax Activation Function - Numerical Stability

- Naive implementation can cause numerical overflow:

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^{n} e^{x_j}}$$

- Stabilized implementation:

$$\text{softmax}(x_i) = \frac{e^{x_i - \max(x)}}{\sum_{j=1}^{n} e^{x_j - \max(x)}}$$

- Subtracting the maximum value prevents overflow:
  - Doesn't change the output (same relative proportions)
  - Ensures at least one exponent equals 1 (when $x_i = \max(x)$)
  - Makes largest exponent manageable
  - Critical for stable computation with larger numbers

# Softmax Activation Function - Gradient

- The gradient is a Jacobian matrix:

$$\frac{\partial \text{softmax}(x_i)}{\partial x_j} = \begin{cases} \text{softmax}(x_i)(1 - \text{softmax}(x_i)) & \text{if } i = j \\ -\text{softmax}(x_i)\text{softmax}(x_j) & \text{if } i \neq j \end{cases}$$

- In matrix form:

$$\nabla_{\mathbf{x}}\text{softmax}(\mathbf{x}) = \text{diag}(\text{softmax}(\mathbf{x})) - \text{softmax}(\mathbf{x}) \otimes \text{softmax}(\mathbf{x})$$

- In practice, often combined with cross-entropy loss for simplification
- When used with cross-entropy, gradient simplifies to: $(\hat{y} - y)$

# Activation Functions - Softmax

```python
class Softmax(Activation):
    """Softmax activation function."""

    def forward(self, x):
        """Forward pass for Softmax function.

        Args:
            x: Input numpy array.
        """
        exp_values = np.exp(x - np.max(x, axis=1, keepdims=True))
        self.output = exp_values / np.sum(exp_values,
                                          axis=1, keepdims=True)
        return self.output

    def backward(self, output_grad):
        """Backward pass for Softmax function.

        Args:
            output_grad: Gradient of the cost with respect to the
        output.
        """
        # Simplified backward pass
        return output_grad
```

# Neural Network Layers - Overview

- Building blocks of neural networks
- Base Layer interface
- Implementations:
  - Dense (Fully connected)
  - Dropout (Regularization)
  - Flatten (Reshaping)
- Each layer implements forward and backward passes

## Layers - Base Class

```python
class Layer(object):
    """Base class for all layers."""

    def __init__(self):
        """Initialize the layer."""
        self.params = {}
        self.grads = {}

    def forward(self, inputs):
        """Forward pass.

        Args:
            inputs: Input data.
        """
        raise NotImplementedError

    def backward(self, output_grad):
        """Backward pass.

        Args:
            output_grad: Gradient of the cost with respect to the
    output.
        """
        raise NotImplementedError
```

# Layers - Dense Layer (1/2)

```python
1  class Dense(Layer):
2      """Fully connected layer."""
3
4      def __init__(self, n_units, input_shape=None,
5                   weight_initializer=None,
6                   bias_initializer=None):
7          """Initialize the dense layer.
8
9          Args:
10             n_units: Number of output units.
11             input_shape: Shape of the input data.
12             weight_initializer: Weight initializer.
13             bias_initializer: Bias initializer.
14         """
15         super(Dense, self).__init__()
16         self.n_units = n_units
17         self.input_shape = input_shape
18
19         self.weight_initializer = weight_initializer
20         if self.weight_initializer is None:
21             self.weight_initializer = HeNormal()
22
23         self.bias_initializer = bias_initializer
24         if self.bias_initializer is None:
25             self.bias_initializer = Zero()
```

# Layers - Dense Layer (2/2)

```python
1    def forward(self, inputs):
2        """Forward pass for dense layer.
3
4        Args:
5            inputs: Input data.
6        """
7        self.inputs = inputs
8
9        if not hasattr(self, 'weights'):
10            self.weights = self.weight_initializer(
11                (self.input_shape, self.n_units))
12            self.bias = self.bias_initializer((1, self.n_units))
13
14            self.params['weights'] = self.weights
15            self.params['bias'] = self.bias
16
17        return np.dot(self.inputs, self.weights) + self.bias
18
19    def backward(self, output_grad):
20        """Backward pass for dense layer.
21
22        Args:
23            output_grad: Gradient of the cost with respect to the
     output.
24        """
```

## Layers - Dropout Layer

```python
class Dropout(Layer):
    """Dropout layer."""

    def __init__(self, dropout_rate):
        """Initialize the dropout layer.

        Args:
            dropout_rate: Dropout rate.
        """
        super(Dropout, self).__init__()
        self.dropout_rate = dropout_rate

    def forward(self, inputs, training=True):
        """Forward pass for dropout layer.

        Args:
            inputs: Input data.
            training: Whether in training mode.
        """
        self.inputs = inputs

        if training:
            self.mask = np.random.binomial(
                1, 1 - self.dropout_rate, size=inputs.shape) / (1 -
    self.dropout_rate)
```

# Loss Functions - Overview

- Measure model performance
- Base Loss abstract class
- Common loss functions:
    - Mean Squared Error (MSE)
    - Categorical Cross-Entropy
    - Binary Cross-Entropy
- Each implements forward and gradient calculations

# Loss Functions - Base Class

```python
class Loss(object):
    """Base class for all loss functions."""

    def forward(self, y_true, y_pred):
        """Forward pass.

        Args:
            y_true: Ground truth values.
            y_pred: Predicted values.
        """
        raise NotImplementedError

    def gradient(self, y_true, y_pred):
        """Gradient of the loss function.

        Args:
            y_true: Ground truth values.
            y_pred: Predicted values.
        """
        raise NotImplementedError
```

# Loss Functions - MSE

```python
1   class MeanSquaredError(Loss):
2       """Mean squared error loss function."""
3
4       def forward(self, y_true, y_pred):
5           """Forward pass for mean squared error.
6
7           Args:
8               y_true: Ground truth values.
9               y_pred: Predicted values.
10          """
11          return 0.5 * np.power(y_pred - y_true, 2).mean()
12
13      def gradient(self, y_true, y_pred):
14          """Gradient of the mean squared error.
15
16          Args:
17              y_true: Ground truth values.
18              y_pred: Predicted values.
19          """
20          return y_pred - y_true
```

# Loss Functions - Categorical Cross-Entropy

```python
1  class CategoricalCrossentropy(Loss):
2      """Categorical cross-entropy loss function."""
3
4      def forward(self, y_true, y_pred):
5          """Forward pass for categorical cross-entropy.
6
7          Args:
8              y_true: Ground truth values.
9              y_pred: Predicted values.
10         """
11         # Clip to avoid log(0)
12         y_pred = np.clip(y_pred, 1e-15, 1 - 1e-15)
13         return -np.sum(y_true * np.log(y_pred)) / y_true.shape[0]
14
15     def gradient(self, y_true, y_pred):
16         """Gradient of the categorical cross-entropy.
17
18         Args:
19             y_true: Ground truth values.
20             y_pred: Predicted values.
21         """
22         # Clip to avoid division by zero
23         y_pred = np.clip(y_pred, 1e-15, 1 - 1e-15)
24         return -y_true / y_pred / y_true.shape[0]
```

# Models - Overview

- Sequential model for chaining layers
- Methods for training and evaluation
- Forward/backward pass implementation
- Training loop with batch processing
- Model evaluation and prediction

# Models - Sequential (1/2)

```python
1  class Sequential(object):
2      """Sequential model."""
3
4      def __init__(self, layers=None):
5          """Initialize the model.
6
7          Args:
8              layers: List of layers.
9          """
10         self.layers = layers if layers is not None else []
11
12     def add(self, layer):
13         """Add a layer to the model.
14
15         Args:
16             layer: Layer to add.
17         """
18         self.layers.append(layer)
```

# Models - Sequential (2/2)

```python
def forward(self, inputs, training=True):
    """Forward pass.

    Args:
        inputs: Input data.
        training: Whether in training mode.
    """
    for layer in self.layers:
        if hasattr(layer, 'training'):
            inputs = layer.forward(inputs, training)
        else:
            inputs = layer.forward(inputs)
    return inputs

def backward(self, grad):
    """Backward pass.

    Args:
        grad: Gradient of the cost with respect to the output.
    """
    for layer in reversed(self.layers):
        grad = layer.backward(grad)
    return grad
```

## Models - Training

```python
def fit(self, x, y, epochs=100, batch_size=32,
        loss_fn=None, optimizer=None,
        validation_data=None, verbose=True):
    """Train the model.

    Args:
        x: Input data.
        y: Target data.
        epochs: Number of epochs.
        batch_size: Batch size.
        loss_fn: Loss function.
        optimizer: Optimizer.
        validation_data: Validation data.
        verbose: Whether to print progress.
    """
    if loss_fn is None:
        loss_fn = MeanSquaredError()

    if optimizer is None:
        optimizer = SGD()

    # Training loop implementation
    for epoch in range(epochs):
        # Process mini-batches
        # Update weights using optimizer
```

# Optimizers - Overview

- Update model parameters based on gradients
- Base Optimizer abstract class
- Common optimization algorithms:
    - SGD (Stochastic Gradient Descent)
    - Adam (Adaptive Moment Estimation)
    - RMSprop
- Each implements the update method

# Optimizers - Base Class

```python
1  class Optimizer(object):
2      """Base class for all optimizers."""
3
4      def __init__(self, learning_rate=0.01):
5          """Initialize the optimizer.
6
7          Args:
8              learning_rate: Learning rate.
9          """
10         self.learning_rate = learning_rate
11
12     def update(self, layer):
13         """Update the layer weights.
14
15         Args:
16             layer: Layer to update.
17         """
18         raise NotImplementedError
```

# Optimizers - SGD

```python
1  class SGD(Optimizer):
2      """Stochastic gradient descent optimizer."""
3
4      def __init__(self, learning_rate=0.01, momentum=0.0):
5          """Initialize the SGD optimizer.
6
7          Args:
8              learning_rate: Learning rate.
9              momentum: Momentum factor.
10         """
11         super(SGD, self).__init__(learning_rate)
12         self.momentum = momentum
13         self.velocity = {}
14
15     def update(self, layer):
16         """Update the layer weights.
17
18         Args:
19             layer: Layer to update.
20         """
21         for param_name in layer.params:
22             # Initialize velocity for the parameter if not exists
23             if param_name not in self.velocity:
24                 self.velocity[param_name] = np.zeros_like(
25                     layer.params[param_name])
```

# Optimizers - Adam

```python
class Adam(Optimizer):
    """Adam optimizer."""

    def __init__(self, learning_rate=0.001, beta_1=0.9,
                 beta_2=0.999, epsilon=1e-8):
        """Initialize the Adam optimizer.

        Args:
            learning_rate: Learning rate.
            beta_1: Exponential decay rate for first moment.
            beta_2: Exponential decay rate for second moment.
            epsilon: Small constant for numerical stability.
        """
        super(Adam, self).__init__(learning_rate)
        self.beta_1 = beta_1
        self.beta_2 = beta_2
        self.epsilon = epsilon
        self.m = {}   # First moment
        self.v = {}   # Second moment
        self.t = 0    # Timestep

    def update(self, layer):
        """Update implementation with moment calculations"""
```

# Building a Complete Neural Network

```python
1  # Import all components
2  from npdl.models import Sequential
3  from npdl.layers import Dense, Dropout
4  from npdl.activations import ReLU, Softmax
5  from npdl.initializers import HeNormal, Zero
6  from npdl.losses import CategoricalCrossentropy
7  from npdl.optimizers import Adam
8
9  # Create a model
10 model = Sequential()
11 model.add(Dense(128, input_shape=784,
12                 weight_initializer=HeNormal(),
13                 bias_initializer=Zero()))
14 model.add(ReLU())
15 model.add(Dropout(0.2))
16 model.add(Dense(64))
17 model.add(ReLU())
18 model.add(Dropout(0.2))
19 model.add(Dense(10))
20 model.add(Softmax())
21
22 # Compile and train
23 model.fit(x_train, y_train, epochs=10, batch_size=32,
24           loss_fn=CategoricalCrossentropy(),
25           optimizer=Adam(learning_rate=0.001),
```

# Summary and Next Steps

- We've covered the complete implementation of:
  - Activation functions
  - Weight initializers
  - Neural network layers
  - Loss functions
  - Model construction
  - Optimizers
- Possible extensions:
  - Convolutional layers
  - Recurrent layers
  - Batch normalization
  - More advanced optimizers
- Practical exercises to implement and test