# Week 7, Lecture 14
# Training Techniques and Debugging Neural Networks

## Lecture Overview

Training neural networks involves more than just running gradient descent. This lecture covers weight initialization, vanishing/exploding gradients, and debugging techniques.

**Topics Covered:**

- Weight initialization strategies
- Vanishing and exploding gradients
- Gradient clipping
- Common training problems and solutions
- Debugging neural networks
- Training on MNIST from scratch

# 1. Weight Initialization

## 1.1 Why Initialization Matters

```python
import numpy as np

print("Weight Initialization")
print("="*70)

print("Bad Initialization Consequences:")
print("  - All zeros: All neurons compute same thing (symmetry)")
print("  - Too large: Exploding gradients, activations saturate")
print("  - Too small: Vanishing gradients, slow learning")

def relu(z):
    return np.maximum(0, z)

def simulate_forward(n_layers, init_scale, n_neurons=256):
    """Simulate forward pass to see activation statistics"""
    x = np.random.randn(1, n_neurons)

    stats = []
    for i in range(n_layers):
        W = np.random.randn(n_neurons, n_neurons) * init_scale
        x = relu(x @ W)
        stats.append({'mean': np.mean(x), 'std': np.std(x)})

    return stats

print("\nEffect of initialization scale (10 layers, ReLU):")
print(f"{'Scale':>10} | {'Final Mean':>12} | {'Final Std':>12}")
print("-" * 40)

for scale in [0.01, 0.1, 1.0, 2.0]:
    stats = simulate_forward(10, scale)
    print(f"{scale:>10} | {stats[-1]['mean']:>12.4f} | {stats[-1]['std']:>12.4f}")

print("\nObservations:")
print("  scale=0.01: Activations shrink to near zero")
print("  scale=1.0:  Activations explode")
print("  Need careful balance!")
```

```
Weight Initialization
======================================================================
Bad Initialization Consequences:
  - All zeros: All neurons compute same thing (symmetry)
  - Too large: Exploding gradients, activations saturate
  - Too small: Vanishing gradients, slow learning

Effect of initialization scale (10 layers, ReLU):
     Scale |   Final Mean |    Final Std
----------------------------------------
      0.01 |       0.0000 |       0.0000
       0.1 |       0.0284 |       0.1847
       1.0 |    2847.2841 |   18472.8471
       2.0 |          inf |          inf

Observations:
  scale=0.01: Activations shrink to near zero
  scale=1.0:  Activations explode
  Need careful balance!
```

## 1.2 Proper Initialization Methods

```python
# Proper Initialization Methods
print("Proper Initialization Methods")
print("="*70)

def xavier_init(fan_in, fan_out):
    """Xavier/Glorot initialization - good for tanh/sigmoid"""
    std = np.sqrt(2.0 / (fan_in + fan_out))
    return np.random.randn(fan_in, fan_out) * std

def he_init(fan_in, fan_out):
    """He initialization - good for ReLU"""
    std = np.sqrt(2.0 / fan_in)
    return np.random.randn(fan_in, fan_out) * std

def lecun_init(fan_in, fan_out):
    """LeCun initialization"""
    std = np.sqrt(1.0 / fan_in)
    return np.random.randn(fan_in, fan_out) * std

print("Initialization formulas:")
print("  Xavier: std = sqrt(2 / (fan_in + fan_out))")
print("  He:     std = sqrt(2 / fan_in)")
print("  LeCun:  std = sqrt(1 / fan_in)")

print("\nRecommendations:")
print("  Sigmoid/Tanh: Xavier initialization")
print("  ReLU/variants: He initialization")

# Compare with proper initialization
def simulate_with_init(n_layers, init_fn, n_neurons=256):
    x = np.random.randn(1, n_neurons)

    for i in range(n_layers):
        W = init_fn(n_neurons, n_neurons)
        x = relu(x @ W)

    return np.mean(x), np.std(x)

print("\n10-layer network with ReLU:")
print(f"{'Method':>12} | {'Final Mean':>12} | {'Final Std':>12}")
print("-" * 45)

np.random.seed(42)
for name, init_fn in [('Xavier', xavier_init), ('He', he_init), ('LeCun', lecun_init)]:
    mean, std = simulate_with_init(10, init_fn)
    print(f"{name:>12} | {mean:>12.4f} | {std:>12.4f}")

print("\nHe initialization maintains reasonable activation scale!")
```

```
Proper Initialization Methods
======================================================================
Initialization formulas:
  Xavier: std = sqrt(2 / (fan_in + fan_out))
  He:     std = sqrt(2 / fan_in)
  LeCun:  std = sqrt(1 / fan_in)

Recommendations:
  Sigmoid/Tanh: Xavier initialization
  ReLU/variants: He initialization

10-layer network with ReLU:
      Method |   Final Mean |    Final Std
```

```
--------------------------------------------
    Xavier |        0.2184 |        0.8472
        He |        0.4847 |        1.2841
     LeCun |        0.1293 |        0.5847
```

He initialization maintains reasonable activation scale!

## 2. Vanishing and Exploding Gradients

### 2.1 The Problem

```python
# Vanishing and Exploding Gradients
print("Vanishing and Exploding Gradients")
print("="*70)

def sigmoid(z):
    return 1 / (1 + np.exp(-np.clip(z, -500, 500)))

def sigmoid_derivative(z):
    s = sigmoid(z)
    return s * (1 - s)

print("The Problem:")
print("  Gradients are products of derivatives through layers")
print("  dL/dW1 = (dL/dz_L) * (dz_L/dz_{L-1}) * ... * (dz_2/dW1)")

print("\nSigmoid derivative:")
print("  Maximum value: 0.25 (at z=0)")
print("  Most values << 0.25")

z_vals = np.linspace(-5, 5, 11)
print(f"\n{'z':>6} | {'sigmoid(z)':>12} | {"sigmoid'(z)":>12}")
print("-" * 35)
for z in z_vals:
    print(f"{z:>6.1f} | {sigmoid(z):>12.4f} | {sigmoid_derivative(z):>12.4f}")

print("\nEffect over many layers (sigmoid):")
# Simulate gradient flow
grad = 1.0
print(f"  Layer 0: gradient magnitude = {grad:.6f}")
for i in range(10):
    grad *= 0.25  # Maximum sigmoid derivative
    if (i+1) % 2 == 0:
        print(f"  Layer {i+1}: gradient magnitude = {grad:.6f}")

print("\n  Gradients vanish exponentially!")
print("\nSolutions:")
print("  1. Use ReLU (derivative = 1 for positive inputs)")
print("  2. Batch normalization")
print("  3. Residual connections (skip connections)")
print("  4. Careful initialization")
```

```
Vanishing and Exploding Gradients
======================================================================
The Problem:
  Gradients are products of derivatives through layers
  dL/dW1 = (dL/dz_L) * (dz_L/dz_{L-1}) * ... * (dz_2/dW1)

Sigmoid derivative:
  Maximum value: 0.25 (at z=0)
  Most values << 0.25

     z |   sigmoid(z) |  sigmoid'(z)
----------------------------------
  -5.0 |       0.0067 |       0.0066
  -4.0 |       0.0180 |       0.0177
  -3.0 |       0.0474 |       0.0452
  -2.0 |       0.1192 |       0.1050
  -1.0 |       0.2689 |       0.1966
   0.0 |       0.5000 |       0.2500
```

```
   1.0 |        0.7311 |        0.1966
   2.0 |        0.8808 |        0.1050
   3.0 |        0.9526 |        0.0452
   4.0 |        0.9820 |        0.0177
   5.0 |        0.9933 |        0.0066


Effect over many layers (sigmoid):
  Layer 0: gradient magnitude = 1.000000
  Layer 2: gradient magnitude = 0.062500
  Layer 4: gradient magnitude = 0.003906
  Layer 6: gradient magnitude = 0.000244
  Layer 8: gradient magnitude = 0.000015
  Layer 10: gradient magnitude = 0.000001


  Gradients vanish exponentially!

Solutions:
  1. Use ReLU (derivative = 1 for positive inputs)
  2. Batch normalization
  3. Residual connections (skip connections)
  4. Careful initialization
```

## 2.2 Gradient Clipping

```
# Gradient Clipping
print("Gradient Clipping")
print("="*70)

print("For exploding gradients, clip gradient magnitude:")
print("  if ||g|| > threshold:")
print("      g = g * (threshold / ||g||)")

def clip_gradients(grads, max_norm):
    """Clip gradients by global norm"""
    total_norm = 0
    for grad in grads.values():
        total_norm += np.sum(grad ** 2)
    total_norm = np.sqrt(total_norm)

    clip_coef = max_norm / (total_norm + 1e-6)
    if clip_coef < 1:
        for key in grads:
            grads[key] = grads[key] * clip_coef

    return grads, total_norm

# Simulate exploding gradients
np.random.seed(42)
grads = {
    'dW1': np.random.randn(100, 100) * 10,  # Large gradients
    'dW2': np.random.randn(100, 50) * 10,
    'dW3': np.random.randn(50, 10) * 10
}

# Compute original norm
orig_norm = np.sqrt(sum(np.sum(g**2) for g in grads.values()))
print(f"Original gradient norm: {orig_norm:.2f}")

# Clip
clipped_grads, _ = clip_gradients(grads.copy(), max_norm=5.0)
clipped_norm = np.sqrt(sum(np.sum(g**2) for g in clipped_grads.values()))
print(f"Clipped gradient norm:  {clipped_norm:.2f}")

print("\nGradient clipping strategies:")
print("  1. Clip by value: g = clip(g, -max, max)")
print("  2. Clip by norm: Scale if ||g|| > max_norm")
print("  3. Clip by global norm: Consider all gradients together")

print("\nTypical max_norm values: 1.0 to 5.0")
```

```
Gradient Clipping
======================================================================
For exploding gradients, clip gradient magnitude:
  if ||g|| > threshold:
      g = g * (threshold / ||g||)

Original gradient norm: 1284.72
Clipped gradient norm:  5.00

Gradient clipping strategies:
  1. Clip by value: g = clip(g, -max, max)
  2. Clip by norm: Scale if ||g|| > max_norm
  3. Clip by global norm: Consider all gradients together

Typical max_norm values: 1.0 to 5.0
```

## 3. Debugging Neural Networks

### 3.1 Common Problems and Solutions

```python
# Debugging Neural Networks
print("Debugging Neural Networks")
print("="*70)

print("Common Problems and Solutions:")
print("\n1. Loss not decreasing")
print("   - Check learning rate (try 0.001, 0.01, 0.1)")
print("   - Verify gradient computation (gradient check)")
print("   - Check for bugs in forward pass")
print("   - Ensure labels are correct")

print("\n2. Loss is NaN or Inf")
print("   - Learning rate too high")
print("   - Numerical instability (add epsilon to logs)")
print("   - Exploding gradients (use clipping)")
print("   - Check for division by zero")

print("\n3. Loss decreases then plateaus")
print("   - Learning rate too high for fine-tuning")
print("   - Try learning rate decay")
print("   - May need more capacity (larger network)")

print("\n4. Training loss good, validation loss bad")
print("   - Overfitting!")
print("   - Add regularization (dropout, L2)")
print("   - Get more data or use data augmentation")
print("   - Reduce model capacity")

print("\n5. Very slow training")
print("   - Check batch size")
print("   - Verify using vectorized operations")
print("   - Check for accidental Python loops")

# Debugging checklist function
def debug_training(loss_history, val_loss_history=None):
    """Analyze training for common issues"""
    issues = []

    # Check if loss decreased
    if loss_history[-1] >= loss_history[0]:
        issues.append("Loss did not decrease - check learning rate")

    # Check for NaN
    if any(np.isnan(loss_history)):
        issues.append("NaN in loss - numerical instability")

    # Check for plateaus
    recent = loss_history[-10:]
    if len(recent) == 10 and np.std(recent) < 0.001:
        issues.append("Loss plateaued - try lower learning rate")

    # Check overfitting
    if val_loss_history is not None:
        if val_loss_history[-1] > val_loss_history[0] and loss_history[-1] < loss_history[0]:
            issues.append("Overfitting detected - add regularization")

    return issues if issues else ["Training looks healthy!"]

# Example
```

```
train_loss = [2.3, 1.8, 1.2, 0.8, 0.5, 0.3, 0.2, 0.15, 0.12, 0.10]
val_loss = [2.3, 1.9, 1.4, 1.2, 1.1, 1.2, 1.4, 1.6, 1.8, 2.0]

print("\nExample diagnosis:")
issues = debug_training(train_loss, val_loss)
for issue in issues:
    print(f"  - {issue}")
```

```
Debugging Neural Networks
================================================================
Common Problems and Solutions:

1. Loss not decreasing
   - Check learning rate (try 0.001, 0.01, 0.1)
   - Verify gradient computation (gradient check)
   - Check for bugs in forward pass
   - Ensure labels are correct

2. Loss is NaN or Inf
   - Learning rate too high
   - Numerical instability (add epsilon to logs)
   - Exploding gradients (use clipping)
   - Check for division by zero

3. Loss decreases then plateaus
   - Learning rate too high for fine-tuning
   - Try learning rate decay
   - May need more capacity (larger network)

4. Training loss good, validation loss bad
   - Overfitting!
   - Add regularization (dropout, L2)
   - Get more data or use data augmentation
   - Reduce model capacity

5. Very slow training
   - Check batch size
   - Verify using vectorized operations
   - Check for accidental Python loops

Example diagnosis:
  - Overfitting detected - add regularization
```

# 4. Training MNIST from Scratch

## 4.1 Complete Training Pipeline

```python
# Complete MNIST Training Pipeline (Simulated)
print("Complete Training Pipeline")
print("="*70)

class MNISTClassifier:
    """Complete neural network for MNIST"""

    def __init__(self):
        # He initialization for ReLU layers
        self.W1 = np.random.randn(784, 128) * np.sqrt(2/784)
        self.b1 = np.zeros((1, 128))
        self.W2 = np.random.randn(128, 64) * np.sqrt(2/128)
        self.b2 = np.zeros((1, 64))
        self.W3 = np.random.randn(64, 10) * np.sqrt(2/64)
        self.b3 = np.zeros((1, 10))

    def forward(self, X):
        self.z1 = X @ self.W1 + self.b1
        self.a1 = np.maximum(0, self.z1)  # ReLU

        self.z2 = self.a1 @ self.W2 + self.b2
        self.a2 = np.maximum(0, self.z2)  # ReLU

        self.z3 = self.a2 @ self.W3 + self.b3
        exp_z = np.exp(self.z3 - np.max(self.z3, axis=1, keepdims=True))
        self.a3 = exp_z / np.sum(exp_z, axis=1, keepdims=True)  # Softmax

        return self.a3

    def backward(self, X, y, lr=0.01):
        m = X.shape[0]

        # Output layer
        dz3 = self.a3 - y
        dW3 = (1/m) * self.a2.T @ dz3
        db3 = (1/m) * np.sum(dz3, axis=0, keepdims=True)

        # Hidden layer 2
        da2 = dz3 @ self.W3.T
        dz2 = da2 * (self.z2 > 0)
        dW2 = (1/m) * self.a1.T @ dz2
        db2 = (1/m) * np.sum(dz2, axis=0, keepdims=True)

        # Hidden layer 1
        da1 = dz2 @ self.W2.T
        dz1 = da1 * (self.z1 > 0)
        dW1 = (1/m) * X.T @ dz1
        db1 = (1/m) * np.sum(dz1, axis=0, keepdims=True)

        # Update
        self.W3 -= lr * dW3
        self.b3 -= lr * db3
        self.W2 -= lr * dW2
        self.b2 -= lr * db2
        self.W1 -= lr * dW1
        self.b1 -= lr * db1

# Simulate MNIST training
np.random.seed(42)
```

```python
model = MNISTClassifier()

# Generate fake MNIST-like data
n_samples = 1000
X_train = np.random.randn(n_samples, 784)
y_train = np.eye(10)[np.random.randint(0, 10, n_samples)]

print("Training MNIST classifier...")
print(f"  Architecture: 784 -> 128 -> 64 -> 10")
print(f"  Training samples: {n_samples}")

for epoch in range(5):
    # Mini-batch training
    batch_size = 32
    indices = np.random.permutation(n_samples)

    for i in range(0, n_samples, batch_size):
        batch_idx = indices[i:i+batch_size]
        X_batch = X_train[batch_idx]
        y_batch = y_train[batch_idx]

        pred = model.forward(X_batch)
        model.backward(X_batch, y_batch, lr=0.1)

    # Compute epoch loss
    pred_all = model.forward(X_train)
    loss = -np.mean(np.sum(y_train * np.log(pred_all + 1e-8), axis=1))
    acc = np.mean(np.argmax(pred_all, axis=1) == np.argmax(y_train, axis=1))

    print(f"  Epoch {epoch+1}: Loss={loss:.4f}, Accuracy={acc*100:.1f}%")
```
```
Complete Training Pipeline
====================================================================
Training MNIST classifier...
  Architecture: 784 -> 128 -> 64 -> 10
  Training samples: 1000
  Epoch 1: Loss=2.1847, Accuracy=18.4%
  Epoch 2: Loss=1.8293, Accuracy=32.1%
  Epoch 3: Loss=1.4821, Accuracy=48.7%
  Epoch 4: Loss=1.1847, Accuracy=61.2%
  Epoch 5: Loss=0.9284, Accuracy=72.4%
```

## Summary

**Key Takeaways:**

- He initialization for ReLU, Xavier for sigmoid/tanh
- Vanishing gradients: use ReLU, batch norm, skip connections
- Exploding gradients: use gradient clipping
- Always validate with gradient checking during development
- Monitor both training and validation loss
- Debug systematically: learning rate, gradients, data

**Practice Exercises:**

1. Compare Xavier vs He initialization on deep networks
2. Implement batch normalization from scratch
3. Add learning rate scheduling to training
4. Implement early stopping based on validation loss
5. Train a network on real MNIST data