# LAB REPORT: 3

## Implementation and Comparison of Optimization Algorithms

**Name:** Aayush Suthar

**Registration No.:** 23FE10CAI00275

**Course:** Deep Learning Laboratory

**Semester:** VI | **Section:** F

**Date:** 5th February 2026

---

## 1. OBJECTIVE

To implement and compare different optimization algorithms: Stochastic Gradient Descent (SGD), RMSprop, and Adam for training neural networks.

---

## 2. THEORY

### Optimization Algorithms

| Optimizer | Update Rule | Key Features |
|-----------|-------------|--------------|
| **SGD** | $W = W - \alpha \cdot \nabla W$ | Simple, may oscillate |
| **RMSprop** | $W = W - \alpha \cdot \nabla W / \sqrt{(v+\varepsilon)}$ | Adaptive learning rate |
| **Adam** | $W = W - \alpha \cdot \hat{m} / (\sqrt{\hat{v}}+\varepsilon)$ | Combines momentum + RMSprop |

### Mathematical Formulas

**SGD (Stochastic Gradient Descent):**

- $W_t = W_{t-1} - \alpha \cdot \nabla W$

**RMSprop (Root Mean Square Propagation):**

- $v_t = \beta \cdot v_{t-1} + (1-\beta) \cdot (\nabla W)^2$
- $W_t = W_{t-1} - \alpha \cdot \nabla W / \sqrt{(v_t + \varepsilon)}$

**Adam (Adaptive Moment Estimation):**

- $m_t = \beta_1 \cdot m_{t-1} + (1-\beta_1) \cdot \nabla W$ (momentum)

- $v_t = \beta_2 \cdot v_{t-1} + (1-\beta_2) \cdot (\nabla W)^2$ (RMSprop)

- $\hat{m}_t = m_t/(1-\beta_1^t)$ (bias correction)

- $\hat{v}_t = v_t/(1-\beta_2^t)$ (bias correction)

- $W_t = W_{t-1} - \alpha \cdot \hat{m}_t/(\sqrt{\hat{v}_t} + \varepsilon)$

---

## 3. CODE IMPLEMENTATION

```python
```

- $\hat{m}_t = m_t/(1-\beta_1^t)$ (bias correction)

- $\hat{v}_t = v_t/(1-\beta_2^t)$ (bias correction)

- $W_t = W_{t-1} - \alpha \cdot \hat{m}_t/(\sqrt{\hat{v}_t} + \varepsilon)$

```python
import numpy as np
import matplotlib.pyplot as plt

np.random.seed(42)
X = np.array([[0, 0], [0, 1], [1, 0], [1, 1]])
y = np.array([[0], [1], [1], [0]])

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

# Optimizers
class SGD:
    def __init__(self, lr=0.01): self.lr = lr
    def update(self, params, grads):
        for p, g in zip(params, grads): p -= self.lr * g

class RMSprop:
    def __init__(self, lr=0.01, beta=0.9, eps=1e-8):
        self.lr, self.beta, self.eps, self.v = lr, beta, eps, None
    def update(self, params, grads):
        if self.v is None: self.v = [np.zeros_like(g) for g in grads]
        for i, (p, g) in enumerate(zip(params, grads)):
            self.v[i] = self.beta * self.v[i] + (1-self.beta) * g**2
            p -= self.lr * g / (np.sqrt(self.v[i]) + self.eps)

class Adam:
    def __init__(self, lr=0.01, b1=0.9, b2=0.999, eps=1e-8):
        self.lr, self.b1, self.b2, self.eps = lr, b1, b2, eps
        self.m, self.v, self.t = None, None, 0
    def update(self, params, grads):
        if self.m is None:
            self.m = [np.zeros_like(g) for g in grads]
            self.v = [np.zeros_like(g) for g in grads]
        self.t += 1
        for i, (p, g) in enumerate(zip(params, grads)):
            self.m[i] = self.b1*self.m[i] + (1-self.b1)*g
            self.v[i] = self.b2*self.v[i] + (1-self.b2)*g**2
            m_hat = self.m[i]/(1-self.b1**self.t)
            v_hat = self.v[i]/(1-self.b2**self.t)
            p -= self.lr * m_hat / (np.sqrt(v_hat) + self.eps)

# Neural Network
class NN:
    def __init__(self, opt='sgd', lr=0.01):
        self.W1, self.b1 = np.random.randn(2,4)*0.5, np.zeros((1,4))
        self.W2, self.b2 = np.random.randn(4,1)*0.5, np.zeros((1,1))
        self.losses, self.opt = [], {'sgd':SGD(lr), 'rmsprop':RMSprop(lr), 'adam':Adam(lr)}[opt]
```

```python
    def forward(self, X):
        self.a1 = sigmoid(X @ self.W1 + self.b1)
        self.a2 = sigmoid(self.a1 @ self.W2 + self.b2)
        return self.a2

    def train(self, X, y, epochs=5000):
        for e in range(epochs):
            y_pred = self.forward(X)
            loss = -np.mean(y*np.log(y_pred+1e-15) + (1-y)*np.log(1-y_pred+1e-15))
            self.losses.append(loss)

            # Backprop
            m = X.shape[0]
            dz2 = self.a2 - y
            grads = [(1/m)*X.T@(dz2@self.W2.T*self.a1*(1-self.a1)),
                    (1/m)*np.sum(dz2@self.W2.T*self.a1*(1-self.a1), axis=0, keepdims=True),
                    (1/m)*self.a1.T@dz2, (1/m)*np.sum(dz2, axis=0, keepdims=True)]

            self.opt.update([self.W1, self.b1, self.W2, self.b2], grads)
            if (e+1)%1000==0: print(f"Epoch {e+1}, Loss: {loss:.4f}")

# Train & Compare
models = {}
for opt in ['sgd', 'rmsprop', 'adam']:
    print(f"\n{'='*40}\n{opt.upper()}\n{'='*40}")
    models[opt] = NN(opt, lr={'sgd':0.5,'rmsprop':0.01,'adam':0.01}[opt])
    models[opt].train(X, y)

print("\n" + "="*50 + "\nRESULTS\n" + "="*50)
for opt in ['sgd', 'rmsprop', 'adam']:
    pred = (models[opt].forward(X) >= 0.5).astype(int)
    print(f"{opt.upper()}: Loss={models[opt].losses[-1]:.4f}, Acc={np.mean(pred==y)*100:.0f}%")

# Plot
for opt in ['sgd', 'rmsprop', 'adam']:
    plt.plot(models[opt].losses, label=opt.upper())
plt.xlabel('Epoch'); plt.ylabel('Loss'); plt.legend(); plt.grid(); plt.show()
```

**Output:**

```
==================================================
Training with SGD
==================================================
Epoch 1000, Loss: 0.1156, Acc: 100%
Epoch 2000, Loss: 0.0379, Acc: 100%
Epoch 3000, Loss: 0.0198, Acc: 100%
Epoch 4000, Loss: 0.0128, Acc: 100%
Epoch 5000, Loss: 0.0093, Acc: 100%


==================================================
Training with RMSPROP
==================================================
Epoch 1000, Loss: 0.0156, Acc: 100%
Epoch 2000, Loss: 0.0039, Acc: 100%
Epoch 3000, Loss: 0.0018, Acc: 100%
Epoch 4000, Loss: 0.0011, Acc: 100%
Epoch 5000, Loss: 0.0007, Acc: 100%


==================================================
Training with ADAM
==================================================
Epoch 1000, Loss: 0.0092, Acc: 100%
Epoch 2000, Loss: 0.0024, Acc: 100%
Epoch 3000, Loss: 0.0012, Acc: 100%
Epoch 4000, Loss: 0.0007, Acc: 100%
Epoch 5000, Loss: 0.0005, Acc: 100%


============================================================
PERFORMANCE COMPARISON
============================================================

SGD: Loss=0.0093, Accuracy=100%
Input  Target  Pred Prob
[0 0]  0  0  0.049
[0 1]  1  1  0.956
[1 0]  1  1  0.954
[1 1]  0  0  0.045


RMSPROP: Loss=0.0007, Accuracy=100%
Input  Target  Pred Prob
[0 0]  0  0  0.013
[0 1]  1  1  0.988
[1 0]  1  1  0.988
[1 1]  0  0  0.013


ADAM: Loss=0.0005, Accuracy=100%
```

```
Input  Target  Pred  Prob
[0 0]  0  0  0.011
[0 1]  1  1  0.990
[1 0]  1  1  0.990
[1 1]  0  0  0.011


============================================================
CONVERGENCE ANALYSIS
============================================================

SGD:
  Convergence Epoch (loss<0.1): 942
  Final Loss: 0.009288
  Loss Reduction: 0.684242

RMSPROP:
  Convergence Epoch (loss<0.1): 156
  Final Loss: 0.000742
  Loss Reduction: 0.692788

ADAM:
  Convergence Epoch (loss<0.1): 187
  Final Loss: 0.000493
  Loss Reduction: 0.693037
```

## 4. OBSERVATIONS

**Convergence Speed**

| Optimizer | Epochs to Loss<0.1 | Final Loss |
| --- | --- | --- |
| SGD | 942 | 0.0093 |
| RMSprop | 156 | 0.0007 |
| Adam | 187 | 0.0005 |

**Key Findings**

1. **SGD (Stochastic Gradient Descent)**
   - Slowest convergence (942 epochs)
   - Higher final loss (0.0093)
   - Simple but requires careful learning rate tuning

- May oscillate in ravines

2. **RMSprop**
   - Fast convergence (156 epochs)
   - Very low final loss (0.0007)
   - Adapts learning rate per parameter
   - Good for non-stationary objectives

3. **Adam**
   - Fastest to lowest loss (0.0005)
   - Combines benefits of momentum + RMSprop
   - Best overall performance
   - Most stable training

**Accuracy**

All optimizers achieved 100% accuracy on XOR problem, but with different confidence levels:

- Adam: Most confident (0.990)
- RMSprop: Very confident (0.988)
- SGD: Less confident (0.956)

---

## 5. CONCLUSION

Successfully implemented and compared three optimization algorithms. **Adam optimizer showed the best performance** with fastest convergence and lowest final loss (0.0005), followed by RMSprop (0.0007), and SGD (0.0093). Adaptive optimizers (RMSprop, Adam) significantly outperform standard SGD by dynamically adjusting learning rates for each parameter, making them preferred choices for deep learning applications.