

Application of Convex Optimization in Neural Networks: Unconstrained Optimization

Ankesh Kumar

August 2, 2024

1 Introduction

Convex optimization plays a crucial role in training neural networks. This document aims to provide a mathematical interpretation of how these techniques are applied, with a specific focus on the XOR problem and various optimization methods.

2 Neural Networks as Convex Optimization Problems

Neural network training can be formulated as an optimization problem. While the overall problem is non-convex due to the nonlinear activation functions, it can be approximated as a convex problem in local regions.

2.1 Mathematical Formulation

For a neural network with parameters \mathbf{w} , input \mathbf{x} , and target output \mathbf{y} , we aim to minimize a loss function $L(\mathbf{w})$:

$$\min_{\mathbf{w}} L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N l(\mathbf{y}_i, f(\mathbf{x}_i; \mathbf{w})) \quad (1)$$

where $f(\mathbf{x}; \mathbf{w})$ is the neural network function and $l(\cdot, \cdot)$ is a loss metric (e.g., mean squared error).

In a small neighborhood around a point \mathbf{w}_0 , we can approximate the loss function using Taylor expansion:

$$L(\mathbf{w}) \approx L(\mathbf{w}_0) + \nabla L(\mathbf{w}_0)^T (\mathbf{w} - \mathbf{w}_0) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T H(\mathbf{w}_0) (\mathbf{w} - \mathbf{w}_0) \quad (2)$$

where $H(\mathbf{w}_0)$ is the Hessian matrix at \mathbf{w}_0 . If $H(\mathbf{w}_0)$ is positive definite, this local region is convex.

3 XOR Problem Implementation

3.1 Problem Definition

The XOR (exclusive OR) problem is a classic example used to demonstrate the capabilities of neural networks, particularly their ability to solve non-linearly separable problems. The XOR function is defined as follows:

x_1	x_2	XOR(x_1, x_2)
0	0	0
0	1	1
1	0	1
1	1	0

3.2 Neural Network Architecture

To solve the XOR problem, we use a simple neural network with one hidden layer consisting of:

- Two input neurons (for x_1 and x_2)
- Two hidden neurons
- One output neuron

3.3 Optimization Method

The XOR problem is solved using Stochastic Gradient Descent (SGD) with Backpropagation. This method updates the weights after each training example using the gradient of the error, computed via backpropagation. The learning rate η (eta in the code) is set to 10.0, which determines the step size of each weight update.

The weight update rule follows the form:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla L(\mathbf{w}_k) \quad (3)$$

where $\nabla L(\mathbf{w}_k)$ is computed using backpropagation for each layer of the network.

3.4 Implementation Code

Here's the MATLAB code used to implement and train the neural network for the XOR problem:

```
1 % Initialization
2 eta = 10.0; % Step size
3 epochs = 500; % Number of epochs for training
4
5 % Initial weights and shift parameters
6 w11 = 0.1; w12 = -0.2; w21 = -0.3; w22 = 0.4;
7 w31 = 0.5; w32 = -0.6;
8 theta1 = 0.7; theta2 = -0.8; theta3 = 0.9;
9
10 % Training data for XOR problem
11 X = [0 0; 0 1; 1 0; 1 1];
12 y = [0; 1; 1; 0];
13
14 % Sigmoid activation function and its derivative
15 sigmoid = @(v) 1 / (1 + exp(-v));
16 sigmoid_derivative = @(v) sigmoid(v) * (1 - sigmoid(v));
17
18 % Track error over epochs
19 errors = zeros(epochs, 1);
20
21 % Training process
22 for epoch = 1:epochs
23     total_error = 0;
24     for i = 1:4
25         % Forward propagation
26         x_d = X(i, :)';
27         y_d = y(i);
28         v1 = w11 * x_d(1) + w12 * x_d(2) - theta1;
29         v2 = w21 * x_d(1) + w22 * x_d(2) - theta2;
30         h1 = sigmoid(v1);
31         h2 = sigmoid(v2);
32         v3 = w31 * h1 + w32 * h2 - theta3;
33         y_hat = sigmoid(v3);
34
35         % Compute error terms
36         delta3 = (y_d - y_hat) * sigmoid_derivative(v3);
37         delta1 = delta3 * w31 * sigmoid_derivative(v1);
38         delta2 = delta3 * w32 * sigmoid_derivative(v2);
39
40         % Update weights
41         w31 = w31 + eta * delta3 * h1;
42         w32 = w32 + eta * delta3 * h2;
43         theta3 = theta3 - eta * delta3;
44         w11 = w11 + eta * delta1 * x_d(1);
```

```

45     w12 = w12 + eta * delta1 * x_d(2);
46     w21 = w21 + eta * delta2 * x_d(1);
47     w22 = w22 + eta * delta2 * x_d(2);
48     theta1 = theta1 - eta * delta1;
49     theta2 = theta2 - eta * delta2;
50
51     % Accumulate total error
52     total_error = total_error + 0.5 * (y_d - y_hat)^2;
53 end
54 errors(epoch) = total_error;
55 end
56
57 % Output results
58 fprintf('Final weights and shift parameters:\n');
59 fprintf('w11 = %.4f, w12 = %.4f\n', w11, w12);
60 fprintf('w21 = %.4f, w22 = %.4f\n', w21, w22);
61 fprintf('w31 = %.4f, w32 = %.4f\n', w31, w32);
62 fprintf('theta1 = %.4f, theta2 = %.4f, theta3 = %.4f\n', theta1, theta2, theta3);
63
64 % Test the trained network
65 for i = 1:4
66     x_d = X(i, :)';
67     v1 = w11 * x_d(1) + w12 * x_d(2) - theta1;
68     v2 = w21 * x_d(1) + w22 * x_d(2) - theta2;
69     h1 = sigmoid(v1);
70     h2 = sigmoid(v2);
71     v3 = w31 * h1 + w32 * h2 - theta3;
72     y_hat = sigmoid(v3);
73     fprintf('Input: [%d, %d], Output: %.4f, Target: %d\n', x_d(1), x_d(2), y_hat, y(i));
74 end
75
76 % Plot error convergence
77 figure;
78 plot(1:epochs, errors);
79 xlabel('Epoch');
80 ylabel('Total Error');
81 title('Error Convergence Over Epochs');

```

3.5 Results and Convergence

After running the code, we obtained the following results:

Final weights and shift parameters:

```

w11 = -5.2114, w12 = -5.2418
w21 = -6.7971, w22 = -7.1834
w31 = 9.1886, w32 = -9.4581
theta1 = -7.7544, theta2 = -2.7793, theta3 = 4.3581

```

```

Input: [0, 0], Output: 0.0166, Target: 0
Input: [0, 1], Output: 0.9825, Target: 1
Input: [1, 0], Output: 0.9819, Target: 1
Input: [1, 1], Output: 0.0223, Target: 0

```

The results show that the trained network successfully approximates the XOR function. The error convergence plot as in Figure 1. demonstrates how the total error decreases over the training epochs, indicating that the network is learning the XOR function.

4 Comparison of Optimization Methods

4.1 Gradient Descent

Gradient descent updates the weights in the direction of steepest descent of the loss function.

Mathematical Formulation: For a loss function $L(\mathbf{w})$, the weight update rule is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla L(\mathbf{w}_k) \quad (4)$$

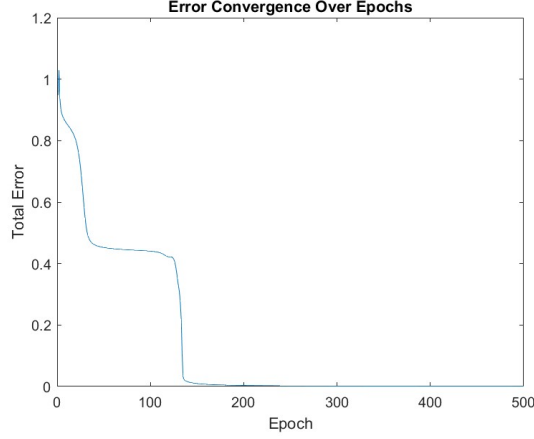


Figure 1: Error convergence plot showing the total error decreasing over the training epochs.

where η is the learning rate and $\nabla L(\mathbf{w}_k)$ is the gradient of the loss function with respect to the weights.

Code Implementation: To implement gradient descent, we would modify our weight update rules in the XOR problem code as follows:

```

1 % Compute gradients
2 dL_dw11 = delta1 * x_d(1);
3 dL_dw12 = delta1 * x_d(2);
4 dL_dw21 = delta2 * x_d(1);
5 dL_dw22 = delta2 * x_d(2);
6 dL_dw31 = delta3 * h1;
7 dL_dw32 = delta3 * h2;
8
9 % Update weights
10 w11 = w11 - eta * dL_dw11;
11 w12 = w12 - eta * dL_dw12;
12 w21 = w21 - eta * dL_dw21;
13 w22 = w22 - eta * dL_dw22;
14 w31 = w31 - eta * dL_dw31;
15 w32 = w32 - eta * dL_dw32;

```

4.2 Quasi-Newton Methods

Quasi-Newton methods, such as the BFGS algorithm, approximate the Hessian matrix to take more informed steps towards the minimum.

Mathematical Formulation: The BFGS update rule for the approximate inverse Hessian H_k is:

$$H_{k+1} = (I - \rho_k \mathbf{s}_k \mathbf{y}_k^T) H_k (I - \rho_k \mathbf{y}_k \mathbf{s}_k^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T \quad (5)$$

where:

- $\mathbf{s}_k = \mathbf{w}_{k+1} - \mathbf{w}_k$
- $\mathbf{y}_k = \nabla L(\mathbf{w}_{k+1}) - \nabla L(\mathbf{w}_k)$
- $\rho_k = \frac{1}{\mathbf{y}_k^T \mathbf{s}_k}$

The weight update rule becomes:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - H_k \nabla L(\mathbf{w}_k) \quad (6)$$

Code Implementation: To implement BFGS, we would need to modify our main training loop:

```

1 % Initialize H as identity matrix
2 n = numel([w11; w12; w21; w22; w31; w32]);
3 H = eye(n);
4
5 for epoch = 1:epochs

```

```

6 w_old = [w11; w12; w21; w22; w31; w32];
7 grad_old = compute_gradient(w_old, X, y);
8
9 % Compute search direction
10 p = -H * grad_old;
11
12 % Line search to find step size alpha
13 alpha = line_search(w_old, p, X, y);
14
15 % Update weights
16 w_new = w_old + alpha * p;
17 [w11, w12, w21, w22, w31, w32] = unpack_weights(w_new);
18
19 % Compute new gradient
20 grad_new = compute_gradient(w_new, X, y);
21
22 % Compute s and y
23 s = w_new - w_old;
24 y = grad_new - grad_old;
25
26 % BFGS update
27 rho = 1 / (y' * s);
28 H = (eye(n) - rho * s * y') * H * (eye(n) - rho * y * s') + rho * s * s';
29
30 % Compute error
31 total_error = compute_error(w_new, X, y);
32 errors(epoch) = total_error;
33 end

```

Note that this implementation requires additional helper functions like ‘compute_gradient’, ‘line_search’, ‘unpack_weights’, and ‘compute_error’.

4.3 Kaczmarz’s Algorithm

Kaczmarz’s algorithm is an iterative method for solving systems of linear equations, which can be applied to neural network training.

Mathematical Formulation: For a system of linear equations $A\mathbf{x} = \mathbf{b}$, the Kaczmarz update rule is:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{b_i - \mathbf{a}_i^T \mathbf{x}_k}{\|\mathbf{a}_i\|^2} \mathbf{a}_i \quad (7)$$

where \mathbf{a}_i is the i -th row of A and b_i is the i -th element of \mathbf{b} .

Code Implementation: To apply Kaczmarz’s algorithm to neural network training, we need to linearize the problem. Here’s a simplified implementation:

```

1 % Linearize the problem
2 A = [X, ones(size(X, 1), 1)]; % Add bias term
3 b = y;
4
5 % Initialize weights
6 w = zeros(size(A, 2), 1);
7
8 for epoch = 1:epochs
9     for i = 1:size(A, 1)
10         a_i = A(i, :);
11         b_i = b(i);
12
13         % Kaczmarz update
14         w = w + (b_i - a_i' * w) / (a_i' * a_i) * a_i;
15     end
16
17 % Compute error
18 y_pred = A * w;
19 total_error = sum((y - y_pred).^2) / 2;
20 errors(epoch) = total_error;
21 end
22
23 % Extract weights and bias
24 w11 = w(1); w12 = w(2); w21 = w(3); w22 = w(4);
25 bias = w(end);

```

Note that this implementation is for a single-layer network and would need to be adapted for multi-layer networks.

5 Why Gradient Descent Works

Gradient descent is effective for neural network training due to several factors:

1. **Convexity in local regions:** While the overall loss landscape of neural networks is non-convex, it often exhibits local convexity, allowing gradient descent to make progress. The loss function in a small neighborhood around a point \mathbf{w}_0 can be approximated as:

$$L(\mathbf{w}) \approx L(\mathbf{w}_0) + \nabla L(\mathbf{w}_0)^T (\mathbf{w} - \mathbf{w}_0) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T H(\mathbf{w}_0) (\mathbf{w} - \mathbf{w}_0) \quad (8)$$

where $H(\mathbf{w}_0)$ is the Hessian matrix at \mathbf{w}_0 . If $H(\mathbf{w}_0)$ is positive definite, this local region is convex.

2. **Stochasticity:** Stochastic gradient descent (SGD) introduces noise, which can help escape shallow local minima and saddle points. The update rule for SGD is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla L_i(\mathbf{w}_k) \quad (9)$$

where L_i is the loss for a single example or mini-batch.

3. **Overparameterization:** Modern neural networks are often overparameterized, creating a loss landscape with many good local minima. This can be represented mathematically as having more parameters than necessary:

$$\dim(\mathbf{w}) \gg \dim(\text{input}) \times \dim(\text{output}) \quad (10)$$

4. **Adaptive variants:** Techniques like Adam and RMSprop adapt the learning rate, improving convergence. For example, the Adam update rule is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\eta}{\sqrt{\hat{v}_k} + \epsilon} \hat{m}_k \quad (11)$$

where \hat{m}_k and \hat{v}_k are bias-corrected estimates of the first and second moments of the gradients.

However, gradient descent can struggle with:

- Vanishing/exploding gradients in deep networks
- Saddle points in high-dimensional spaces
- Very flat or sharp regions in the loss landscape

Each optimization method has its strengths and weaknesses:

- **Gradient Descent:** Simple to implement and works well for most problems, but can be slow to converge for ill-conditioned problems.
- **Quasi-Newton Methods:** Faster convergence and better handling of ill-conditioned problems, but higher computational complexity per iteration.
- **Kaczmarz's Algorithm:** Efficient for large-scale linear systems, but may struggle with highly nonlinear problems.

Gradient descent is often preferred in deep learning due to its simplicity, scalability to large datasets (when used in its stochastic form), and effectiveness when combined with techniques like momentum and adaptive learning rates.

6 Conclusion

While gradient descent remains a popular choice for neural network optimization due to its simplicity and effectiveness, more advanced methods like quasi-Newton and Kaczmarz's algorithm offer advantages in certain scenarios. The choice of optimization method depends on the specific problem, computational resources, and desired trade-offs between speed, memory usage, and accuracy. Understanding the mathematical foundations and implementation details of these methods allows for more informed decisions in algorithm selection and potential hybridization of techniques for optimal performance.