

MODULE-II

TRAINING DEEP NEURAL NETWORKS

Backpropagation with Computational Graph:

Computational graphs are a fundamental tool in deep learning and machine learning, essential for representing and managing complex computations. They play a crucial role in understanding and implementing backpropagation, which is vital for training neural networks. Here's an overview of computational graphs and their significance in backpropagation:

What is a Computational Graph?

A computational graph is a directed graph where:

- **Nodes** represent operations (such as addition or multiplication) or variables (such as inputs or weights).
- **Edges** represent the flow of data (tensors) between these operations.

This graph captures the sequence of operations and their dependencies, allowing for a clear visualization and computation of both forward and backward passes in a neural network.

Why Use Computational Graphs?

1. Organize Complex Computations:

- Neural networks involve numerous layers and operations, making them complex to manage. Computational graphs offer a structured approach to break down these computations into manageable components.

2. Efficient Gradient Computation:

- Backpropagation depends on the efficient calculation of gradients. Computational graphs facilitate automatic differentiation by systematically applying the chain rule across the graph to compute gradients for all parameters.

3. Modularity and Reusability:

- Graphs allow for modular model construction. Different parts of the model can be represented as subgraphs and reused, making it easier to experiment with and build models.

4. Optimization:

- Computational graphs can be optimized to enhance performance. For example, common subexpressions can be shared, and computations can be tailored for specific hardware (such as GPUs).

5. Debugging and Visualization:

- Graphs provide a visual representation of the computation process. This helps in understanding the model, debugging issues, and ensuring that both the forward and backward passes are correctly implemented.

How Computational Graphs Work

1. Forward Pass:

- **Description:** Computes the network's output by propagating inputs through the graph, performing the operations defined by the nodes, and following the edges.
- **Implementation:** Evaluates each node according to the operation it represents, using inputs and intermediate results from connected nodes.

2. Backward Pass (Backpropagation):

- **Description:** Computes the gradients of the loss with respect to each parameter by propagating gradients backward through the graph. This involves applying the chain rule of calculus to the operations represented by the nodes.
- **Implementation:** Traverses the graph in reverse order, starting from the loss node and moving toward the input nodes, updating gradients at each node based on the gradients from subsequent nodes.

Advantages

1. Automatic Differentiation:

- Computational graphs support automatic differentiation, which is crucial for efficiently computing gradients required for training neural networks.

2. Flexibility:

- They allow for dynamic model definitions and modifications, making them suitable for complex and variable neural network architectures.

3. Scalability:

- Optimizations and parallelization strategies can be applied to the graph to enhance computational efficiency, particularly for large models and datasets.

Backpropagation with Post-Activation and Pre-Activation Variables

In neural networks, backpropagation involves computing gradients to update the model's parameters. The process can vary depending on whether gradients are computed with respect to

post-activation or pre-activation variables. Understanding these concepts is crucial for the efficient implementation and optimization of neural networks.

Key Concepts

1. Pre-Activation:

- The raw output of a linear transformation before applying the activation function.
- For a given layer l :

$$z^l = W^l a^{l-1} + b^l$$

where W^l is the weight matrix, a^{l-1} is the activation from the previous layer, and b^l is the bias.

- The pre-activation z^l is then passed through an activation function σ to produce the post-activation:

$$a^l = \sigma(z^l)$$

2. Post-Activation:

- The output of the activation function applied to the pre-activation.
- For the same layer l :

$$\downarrow i^l = \sigma(z^l)$$

Backpropagation with Pre-Activation and Post-Activation Variables

Backpropagation involves computing gradients with respect to each layer's parameters and inputs. This can be done either using pre-activation or post-activation variables.

Backpropagation with Post-Activation Variables

1. Forward Pass:

- Compute pre-activation: $z^l = W^l a^{l-1} + b^l$
- Apply activation function: $a^l = \sigma(z^l)$

2. Gradient Computation:

- The gradient of the loss L with respect to post-activation a^l is needed to propagate errors backward:

$$\frac{\partial L}{\partial a^l}$$

- Compute gradient of the activation function σ :

$$\frac{\partial a^l}{\partial z^l} = \sigma'(z^l)$$

- Calculate gradient of the loss with respect to the pre-activation z^l :

$$\frac{\partial L}{\partial z^l} = \frac{\partial L}{\partial a^l} \cdot \sigma'(z^l)$$

3. Parameter Updates:

- Compute gradients with respect to weights and biases:

$$\frac{\partial L}{\partial W^l} = \frac{\partial L}{\partial z^l} \cdot (a^{l-1})^T$$

$$\frac{\partial L}{\partial b^l} = \frac{\partial L}{\partial z^l}$$

- Update weights and biases using gradient descent:

$$W^l \leftarrow W^l - \eta \frac{\partial L}{\partial W^l}$$

$$b^l \leftarrow b^l - \eta \frac{\partial L}{\partial b^l}$$

Gradient based Strategies:

Learning Rate Decay

Learning rate decay is a technique used to adjust the learning rate during neural network training. It involves gradually reducing the learning rate as training progresses. The strategy is to start with a relatively high learning rate to accelerate convergence and then decrease it for more precise adjustments as the model nears the optimal solution.

Why Use Learning Rate Decay?

- Initial Fast Learning: A larger learning rate at the beginning of training allows the optimizer to make substantial updates quickly, which speeds up learning and helps avoid shallow local minima.
- Later Precise Adjustments: As the model approaches the minimum, a smaller learning rate enables finer adjustments, reducing the risk of overshooting and promoting a smoother convergence.

Without decay, a high learning rate might cause the model to oscillate around the minimum, while starting with a low learning rate could make training excessively slow.

Common Learning Rate Decay Strategies

1. Step Decay:

- The learning rate is reduced by a factor at predefined steps (e.g., every few epochs).

- Formula:

$$\text{new_lr} = \text{initial_lr} \times \text{decay_factor}^{\left\lfloor \frac{\text{epoch}}{\text{step_size}} \right\rfloor}$$

- Example: Reduce learning rate by half every 10 epochs.

This leads to a slower reduction compared to exponential decay.

- The learning rate first increases and then decreases over the course of training. This strategy can help speed up convergence and improve generalization.

Example

Let's say you're training a neural network for image classification using **gradient descent**. You want to use step decay to adjust the learning rate as the training progresses. You decide to reduce the learning rate by a factor of 0.5 every 10 epochs.

- **Initial Learning Rate (LR):** 0.1
- **Decay Factor:** 0.5
- **Step Size (epochs):** 10

Example Values:

Let's calculate the learning rate at various epochs.

1. Epoch 0 to 9:

- $\left\lfloor \frac{t}{10} \right\rfloor = 0$
- Learning Rate: $0.1 \times 0.5^0 = 0.1$

2. Epoch 10 to 19:

- $\left\lfloor \frac{t}{10} \right\rfloor = 1$
- Learning Rate: $0.1 \times 0.5^1 = 0.05$

3. Epoch 20 to 29:

- $\left\lfloor \frac{t}{10} \right\rfloor = 2$
- Learning Rate: $0.1 \times 0.5^2 = 0.025$

4. Epoch 30 to 39:

- $\left\lfloor \frac{t}{10} \right\rfloor = 3$



Momentum-based learning:

Momentum-based learning is an optimization technique used to enhance the performance and convergence of training deep learning models. It refines the gradient descent algorithm by incorporating past gradient information into the current update.

Concept of Momentum-Based Learning

Momentum-based learning involves adding a fraction of the previous gradient update to the current gradient update. This approach smooths the optimization path, reduces oscillations, and accelerates convergence, particularly in areas with small gradients or narrow valleys.

Mathematical Formulation

The update rules for momentum-based learning are as follows:

1. Velocity Update:

- Compute the velocity (a running average of past gradients):

$$v_t = \beta v_{t-1} + (1 - \beta) \nabla L$$

- Here, v_t is the velocity at time t , β is the momentum factor, ∇L is the gradient of the loss function, and v_{t-1} is the velocity from the previous iteration.

2. Parameter Update:

- Update the model parameters using the velocity:

$$\theta_{t+1} = \theta_t - \alpha v_t$$

- Here, θ_{t+1} is the updated parameter, θ_t is the current parameter, and α is the learning rate.

Key Components

1. Momentum Factor (β):

- The momentum factor (β) controls the influence of previous gradients on the current update. It typically ranges from 0.5 to 0.99.
- A higher value of β gives more weight to past gradients, making the updates smoother but potentially slower to react to recent changes.

2. Learning Rate (α):

- The learning rate controls the step size for each update. Momentum helps to navigate smoother paths but still requires careful tuning of the learning rate to avoid overshooting.

Benefits of Momentum-Based Learning

Accelerated Convergence:

Momentum helps accelerate gradient vectors in the correct direction, leading to faster convergence, especially in regions with small gradients or along long, narrow valleys.

Reduced

Oscillations:

By incorporating past gradients, momentum-based learning minimizes oscillations and smooths the path towards the minimum.

Improved Performance:
Momentum aids in escaping local minima and navigating complex loss landscapes more effectively.

Parameter-Specific Learning Rates (PSLR)

The learning rate is a critical hyperparameter that controls the extent of weight adjustments based on the gradient of the loss function. A smaller learning rate slows learning but can lead to precise convergence, while a larger learning rate speeds up learning but may overshoot optimal solutions.

Traditionally, a single global learning rate is used for all model parameters. However, different parts of a neural network may benefit from distinct learning rates. This is where Parameter-Specific Learning Rates (PSLR) come into play, involving assigning different learning rates to various subsets of parameters.

Why Use Parameter-Specific Learning Rates?

Fine-Tuning **Pretrained** **Models:**
When fine-tuning a pretrained model, lower layers (with general features) are often frozen, while higher layers (specific to the new task) are adjusted. Using lower learning rates for pretrained layers prevents disrupting learned features, while higher learning rates can be applied to newly added layers.

Different Learning Rates for Different Layers:
In deep networks, lower layers represent general features and may require smaller learning rates for stability, whereas higher layers, being more task-specific, might benefit from larger learning rates for quicker adaptation.

Layer Types with Different Dynamics:
Convolutional, recurrent, and fully connected layers have different roles and may require different learning rates. For example, recurrent layers like LSTMs or GRUs may need smaller learning rates to handle sensitive updates compared to convolutional or fully connected layers.

Regularization:

Combining parameter-specific learning rates with parameter-specific regularization can provide

optimized and controlled learning dynamics, particularly when different parts of the network have varying levels of regularization.

Complex Architectures:

In complex architectures like ResNets, GANs, transformers, or multi-branch networks, different branches or components may need specific learning rates. PSLR helps manage the learning dynamics of these complex interactions more effectively.

How Parameter-Specific Learning Rates Work

Parameter-specific learning rates involve grouping parameters and assigning each group its own learning rate. This approach allows different parameter groups to be updated at different rates during training.

Example Implementation:

Layer 1: Learning rate of 0.001

Layer 2: Learning rate of 0.01

In this setup, Layer 1 (possibly pretrained) is updated slowly, while Layer 2 adapts more rapidly, allowing for tailored training dynamics.

Mathematical Formulation

For parameter-specific learning rates, the update rule for the parameters θ_i is:

$$\theta_{i,t+1} = \theta_{i,t} - \alpha_i \nabla L(\theta_{i,t})$$

- $\theta_{i,t}$: The value of parameter i at iteration t .
- α_i : Learning rate for parameter i .
- $\nabla L(\theta_{i,t})$: Gradient of the loss function with respect to parameter i at iteration t .

Comparison with Global Learning Rate

Aspect	Global Learning Rate	Parameter-Specific Learning Rate
Learning Rate	Single learning rate for all parameters	Different learning rates for different parameters
Flexibility	Less flexible; same rate applied uniformly	More flexible; adapts rates based on parameter needs
Convergence	May be slower if parameters have different sensitivities	Can converge faster by optimizing learning rates for each parameter
Complexity	Simple to implement	More complex to manage and tune
Handling Different Scales	May require manual adjustment	Automatically adapts to different parameter scales

Gradient Clipping

Gradient clipping is a technique used to address the problem of exploding gradients during the training of deep neural networks. Exploding gradients occur when gradients become excessively large during backpropagation, leading to instability in the training process. This can result in excessively large updates to the model parameters, causing divergence and training failure.

Purpose:

The goal of gradient clipping is to cap the gradients to a predefined range or threshold to ensure they do not exceed a certain size. This stabilizes and controls the training process, particularly in deep or recurrent neural networks (RNNs), where exploding gradients are more common.

Types of Gradient Clipping:

1. Gradient Clipping by Value

In this approach, each component of the gradient vector is clipped if it exceeds a specified range, typically defined by a minimum and maximum value. The idea is to limit the magnitude of each gradient element to prevent any from becoming too large.

How it Works:

- Define a minimum threshold (`clip_value_min`) and a maximum threshold (`clip_value_max`).
- If any gradient component exceeds `clip_value_max`, it is set to `clip_value_max`.
- Conversely, if any gradient component is below `clip_value_min`, it is set to `clip_value_min`.

Example:

- Let's assume a gradient vector $\mathbf{g} = [0.5, -1.2, 3.0, -4.5, 2.0]$.
- If we clip the values to lie between -2 and 2, the clipped gradient vector will become $\mathbf{g} = [0.5, -1.2, 2.0, -2.0, 2.0]$.

Formula:

For each gradient element g_i , the clipping formula is:

$$g_i = \max(\min(g_i, \text{clip_value_max}), \text{clip_value_min})$$

- **Gradient Component (gi):** Each individual element of the gradient vector, which corresponds to the gradient with respect to a specific parameter in the model. This clipping is applied separately to each component.
- **clip_value_max:** The upper bound for gradient values. If any gradient component exceeds this maximum value, it is clipped to `clip_value_max`.
- **clip_value_min:** The lower bound for gradient values. If any gradient component is below
- **clip_value_min:** This is the lower bound for the gradient value. Any gradient component smaller than this will be clipped to this minimum value.

Step-by-Step Process:

1. First, take the minimum of the gradient component g_i and the upper bound `clip_value_max`. This ensures that g_i does not exceed the upper bound.
2. Then, take the maximum of the result and the lower bound `clip_value_min`. This ensures that g_i does not go below the lower bound.
3. The result is a gradient value g_i that is within the specified range $[clip_value_min, clip_value_max]$.

Example:

Suppose $g_i = 3.0$, `clip_value_max = 2.0`, and `clip_value_min = -2.0`.

- The first step is to compute $\min(3.0, 2.0)$, which is 2.0.
- Then, compute $\max(2.0, -2.0)$, which remains 2.0.

So, the clipped gradient value is 2.0.

Gradient Clipping by Norm

Gradient Clipping by Norm is a technique used to stabilize training by rescaling the entire gradient vector if its magnitude exceeds a predefined threshold.

How It Works:

1. Compute the Norm: Calculate the norm (typically the Euclidean or L2 norm) of the gradient vector.
2. Compare with Threshold: If the norm of the gradient vector exceeds the threshold (`clip_norm`), rescale the entire vector.
3. Rescaling: Adjust the gradient vector so that its norm equals `clip_norm`, while maintaining its direction.

Advantages:

- Preserves Direction: The gradient's direction is maintained, which helps in preserving the gradient's relative scales and can lead to more stable convergence.
- Effective for Deep Models: Particularly useful for models where gradients can vary greatly in magnitude, such as RNNs or very deep networks.

Disadvantages:

- Complexity: Slightly more complex to implement than clipping by value, as it involves computing norms and rescaling vectors.
- Potential for Inefficiency: If norms are frequently close to or exceed the threshold, it may lead to rescaling more often, which could affect training dynamics.

Example:

- Let's assume the gradient vector `g = [2.0, -4.0, 1.0, -3.0]`.
- The L2 norm of `g` is $\|g\|_2 = \sqrt{2.0^2 + (-4.0)^2 + 1.0^2 + (-3.0)^2} = \sqrt{4 + 16 + 1 + 9} = 5.477$.
- If we set the `clip_norm` to 3.0, the gradient vector will be rescaled by a factor of $\frac{3.0}{5.477} \approx 0.548$.
- The clipped gradient vector will become `g = [1.096, -2.192, 0.548, -1.644]`.

Formula:

Let $\|g\|$ represent the norm of the gradient vector g . The formula for clipping is:

$$g = \frac{g}{\|g\|} \cdot \text{clip_norm} \quad \text{if } \|g\| > \text{clip_norm}$$

Otherwise, the gradient remains unchanged.

Scenario:

Suppose you have a gradient vector g for your model's parameters as follows:

$$g = [3.0, 4.0]$$

This is a 2D gradient vector (for simplicity).

You want to apply gradient clipping by norm with a maximum allowable norm `clip_norm` = 5.0.

Step 1: Calculate the norm of the gradient vector

The norm of the gradient vector g , using the L2 norm (Euclidean norm), is calculated as:

$$\|g\|_2 = \sqrt{3.0^2 + 4.0^2} = \sqrt{9 + 16} = \sqrt{25} = 5.0$$

Step 2: Check if clipping is needed

In this case, the norm of the gradient $\|g\|_2 = 5.0$, which is exactly equal to the clipping threshold `clip_norm` = 5.0.

Since the norm of the gradient does **not exceed** the clipping threshold, **no clipping** is necessary.

The clipped gradient remains the same as the original gradient:

`gclipped=[3.0,4.0]`

New Scenario (where clipping is needed):

Now, let's assume the gradient vector was larger:

$$g = [6.0, 8.0]$$

You still want to apply gradient clipping by norm with `clip_norm` = 5.0.

Step 1: Calculate the norm of the new gradient vector

The L2 norm of this gradient vector is:

$$\|g\|_2 = \sqrt{6.0^2 + 8.0^2} = \sqrt{36 + 64} = \sqrt{100} = 10.0$$

Step 2: Check if clipping is needed

Here, the norm $\|g\|_2 = 10.0$ exceeds the clipping threshold `clip_norm` = 5.0. Thus, clipping is needed.

Step 3: Apply gradient clipping by norm

To clip the gradient, we scale down the entire gradient vector such that its norm becomes equal to the threshold `clip_norm` = 5.0.

The scaling factor is computed as:

$$\text{scaling factor} = \frac{\text{clip_norm}}{\|g\|_2} = \frac{5.0}{10.0} = 0.5$$

Step 4: Scale the gradient

Multiply each component of the gradient vector by the scaling factor:

$$g_{\text{clipped}} = [6.0 \times 0.5, 8.0 \times 0.5] = [3.0, 4.0]$$

Final Result:

After clipping, the new gradient vector is:

$$g_{\text{clipped}} = [3.0, 4.0]$$

This vector has a norm of exactly 5.0, matching the `clip_norm`, and its direction is preserved.

Summary:

- Original gradient vector: [6.0, 8.0]
- Norm of original gradient: 10.0
- Clipped gradient vector: [3.0, 4.0]



Polyak Averaging

Polyak Averaging is a method employed to enhance the convergence and stability of optimization algorithms, especially in iterative methods such as gradient descent. This technique involves computing and using the average of model parameters over multiple iterations to obtain a more stable and potentially improved solution.

Concept of Polyak Averaging:

1. Definition:
 - Polyak Averaging entails calculating an average of the model's parameters across various iterations. Instead of relying solely on the parameters from the final iteration, this average is used for making predictions or as the final model.

2. Mathematical Formulation:

- Suppose θ_t is the parameter vector at iteration t . The Polyak Averaged parameter vector $\bar{\theta}_T$ after T iterations is computed as:

$$\bar{\theta}_T = \frac{1}{T} \sum_{t=1}^T \theta_t$$

- Alternatively, Polyak Averaging can use a weighted average, where more recent iterations have higher weights.

Benefits of Polyak Averaging

1. Improves Stability:

- By averaging parameters over multiple iterations, the effects of noisy updates or fluctuations are mitigated, leading to more stable solutions.

2. Reduces Variance:

- This technique helps smooth out noise in the optimization process, potentially resulting in a more accurate and reliable final model.

Applications of Polyak Averaging

1. Stochastic Gradient Descent (SGD):

- Frequently employed with SGD to enhance the stability of the final outcome.

2. Neural Networks:

- Utilized to improve the performance and stability of models trained through iterative optimization techniques.

Local and Spurious Minima

Local Minima

1. Definition:

- A local minimum is a point in the loss landscape where the loss function has a lower value than at nearby points but is not necessarily the lowest value across the entire landscape.

2. Characteristics:

- **Surrounding Area:** The loss value is lower compared to neighboring points, though it may not be the global minimum.

- **Convergence:** Optimization algorithms may become trapped in local minima without mechanisms to escape.

Spurious Minima

1. Definition:

- A spurious minimum is a local minimum that does not accurately reflect the global structure of the loss landscape, often arising from noise or irregularities in the loss surface.

2. Characteristics:

- **Non-Optimal:** Spurious minima are generally not true solutions and may not lead to models with good generalization.
- **Irregular Landscapes:** Common in complex loss landscapes with significant irregularities.

