

Autograd, Computational Graphs and Vanishing Gradients

A Comprehensive Overview

Sayan Chaki

March 20, 2025

Table of Contents

- 1 Introduction to Automatic Differentiation
- 2 Computational Graphs
- 3 Autograd Implementation
- 4 Vanishing Gradient Problem
- 5 Solutions to Vanishing Gradients
- 6 Practical Examples

What is Automatic Differentiation?

- **Automatic differentiation** (autograd) is a set of techniques to efficiently compute exact derivatives
- Essential component of modern deep learning frameworks
- Allows for training of complex neural networks via gradient-based optimization
- Different from numerical differentiation (finite differences) and symbolic differentiation

Types of Automatic Differentiation

Forward Mode

- Computes derivatives alongside function evaluation
- Efficient for functions with many outputs and few inputs

Reverse Mode

- Computes derivatives after function evaluation
- Efficient for functions with many inputs and few outputs
- Used in most deep learning frameworks

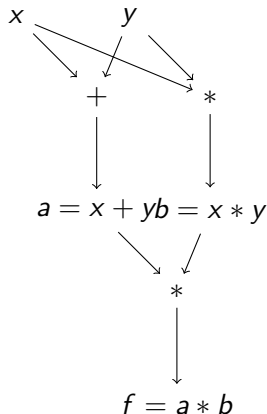
Backpropagation = Reverse-mode autodiff + Chain rule

Computational Graphs

- A computational graph is a directed graph representing mathematical operations
- Nodes: variables (inputs, parameters, intermediate values, outputs)
- Edges: operations that transform values
- Provides a visual and conceptual framework for understanding calculation flow
- Forms the basis for implementing automatic differentiation

Computational Graph Example

Consider the function: $f(x, y) = (x + y) * (x * y)$



Forward Pass vs. Backward Pass

Forward Pass

- Compute values from inputs to outputs
- Store intermediate results for backward pass
- Example:
 $f(x, y) = (x + y) * (x * y)$

$$a = x + y$$

$$b = x * y$$

$$f = a * b$$

Backward Pass

- Compute gradients from outputs to inputs
- Apply chain rule to propagate derivatives

$$\frac{\partial f}{\partial a} = b$$

$$\frac{\partial f}{\partial b} = a$$

$$\begin{aligned}\frac{\partial f}{\partial x} &= \frac{\partial f}{\partial a} \frac{\partial a}{\partial x} + \frac{\partial f}{\partial b} \frac{\partial b}{\partial x} \\ &= b \cdot 1 + a \cdot y\end{aligned}$$

PyTorch Autograd Basics

```
1 import torch
2
3 # Create tensors with autograd enabled
4 x = torch.tensor(2.0, requires_grad=True)
5 y = torch.tensor(3.0, requires_grad=True)
6
7 # Forward pass
8 a = x + y
9 b = x * y
10 f = a * b
11
12 # Compute gradients
13 f.backward()
14
15 # Display results
16 print(f"f(x,y) = {f.item()}")
17 print(f"df/dx = {x.grad.item()}")
18 print(f"df/dy = {y.grad.item()}")
```

Output:

Building a Simple Autograd Engine

```
1 class Value:
2     def __init__(self, data, _children=(), _op=''):
3         self.data = data
4         self.grad = 0
5         self._backward = lambda: None
6         self._prev = set(_children)
7         self._op = _op
8
9     def __add__(self, other):
10        other = other if isinstance(other, Value) else Value
11        (other)
12        out = Value(self.data + other.data, (self, other), '
13        +')
14
15    def _backward():
16        self.grad += out.grad
17        other.grad += out.grad
18        out._backward = _backward
19
20    return out
```

Custom Autograd Example (continued)

```
1  def __mul__(self, other):
2      other = other if isinstance(other, Value) else Value
          (other)
3      out = Value(self.data * other.data, (self, other),
          *')
4
5      def _backward():
6          self.grad += other.data * out.grad
7          other.grad += self.data * out.grad
8      out._backward = _backward
9
10     return out
11
12     def backward(self):
13         # Topological sort
14         topo = []
15         visited = set()
16         def build_topo(v):
17             if v not in visited:
18                 visited.add(v)
19                 for child in v.prev:
```

The Chain Rule in Autograd

The chain rule is fundamental to autograd:

$$\frac{df}{dx} = \frac{df}{dz} \cdot \frac{dz}{dy} \cdot \frac{dy}{dx} \quad (1)$$

For a computational graph:

$$\frac{\partial L}{\partial w_i} = \sum_j \frac{\partial L}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_i} \quad (2)$$

where:

- L is the loss function
- z_j are intermediate nodes
- w_i are parameters

The Vanishing Gradient Problem

- Problem: gradients decrease exponentially as they propagate through layers in deep networks
- Result: earlier layers learn very slowly or not at all
- Causes:
 - Activation functions with small derivatives (e.g., sigmoid, tanh)
 - Weight initialization schemes
 - Deep network architectures
- Consequence: training stalls and performance plateaus

Visualizing Vanishing Gradients

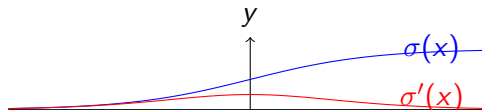
Sigmoid activation function:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (3)$$

Its derivative:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \quad (4)$$

Maximum value of $\sigma'(x)$ is 0.25 at $x = 0$



Gradient Flow in Deep Networks

For a deep network with L layers using sigmoid activations:

$$\left\| \frac{\partial L}{\partial w^{(1)}} \right\| \approx \prod_{l=2}^L \|w^{(l)}\| \cdot \prod_{l=1}^L \|\sigma'(z^{(l)})\| \cdot \left\| \frac{\partial L}{\partial y} \right\| \quad (5)$$

If each $\|\sigma'(z^{(l)})\| \leq 0.25$ and L is large:

$$\prod_{l=1}^L \|\sigma'(z^{(l)})\| \leq 0.25^L \quad (6)$$

For $L = 10$: $0.25^{10} \approx 10^{-6}$ (gradient virtually disappears)

Demonstrating Vanishing Gradients with Code

```
1 import torch
2 import torch.nn as nn
3 import matplotlib.pyplot as plt
4
5 # Create a deep network with sigmoid activations
6 def create_deep_network(depth, activation=nn.Sigmoid()):
7     layers = []
8     for i in range(depth):
9         layers.append(nn.Linear(1, 1))
10        if i < depth - 1:
11            layers.append(activation)
12    return nn.Sequential(*layers)
13
14 # Track gradients through backpropagation
15 def check_gradients(model, depth):
16     x = torch.ones(1, 1, requires_grad=True)
17     y = model(x)
18     y.backward()
19
20     gradients = []
21     for i, layer in enumerate(model):
```

Solutions to the Vanishing Gradient Problem

- **Better activation functions**

- ReLU: $f(x) = \max(0, x)$
- Leaky ReLU: $f(x) = \max(0.01x, x)$
- ELU, SELU, Swish

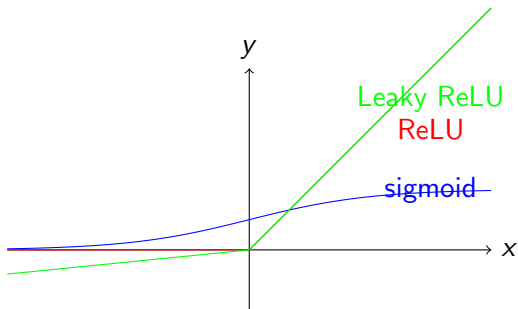
- **Better initialization schemes**

- Xavier/Glorot initialization
- He initialization

- **Architectural innovations**

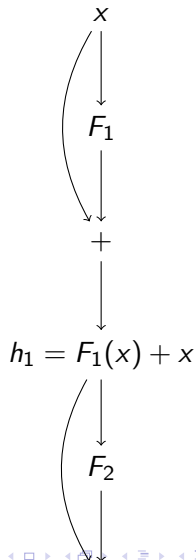
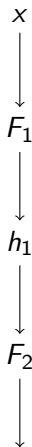
- Skip connections (ResNet)
- Normalization layers (BatchNorm, LayerNorm)
- LSTM/GRU for sequence data

Activation Functions Comparison



Residual Network

Standard Network



Implementing Solutions in PyTorch

```
1 # Better activation function
2 model_relu = nn.Sequential(
3     nn.Linear(784, 256),
4     nn.ReLU(),
5     nn.Linear(256, 128),
6     nn.ReLU(),
7     nn.Linear(128, 10)
8 )
9
10 # Batch normalization
11 model_with_bn = nn.Sequential(
12     nn.Linear(784, 256),
13     nn.BatchNorm1d(256),
14     nn.ReLU(),
15     nn.Linear(256, 128),
16     nn.BatchNorm1d(128),
17     nn.ReLU(),
18     nn.Linear(128, 10)
19 )
20
21 # Residual connections
```

Training a Deep Network with PyTorch

```
1 # Define a fully-connected network with ReLU activations
2 class DeepMLP(nn.Module):
3     def __init__(self, activation=nn.ReLU()):
4         super().__init__()
5         self.model = nn.Sequential(
6             nn.Linear(784, 512),
7             activation,
8             nn.Linear(512, 256),
9             activation,
10            nn.Linear(256, 128),
11            activation,
12            nn.Linear(128, 64),
13            activation,
14            nn.Linear(64, 10)
15        )
16
17    def forward(self, x):
18        return self.model(x)
19
20 # Train with different activation functions
21 relu_model = DeepMLP(activation=nn.ReLU())
```

Visualizing Gradient Flow

```
1 def plot_grad_flow(named_parameters):
2     ave_grads = []
3     layers = []
4     for n, p in named_parameters:
5         if (p.requires_grad) and ("bias" not in n):
6             layers.append(n)
7             ave_grads.append(p.grad.abs().mean().item())
8     plt.figure(figsize=(10, 8))
9     plt.bar(range(len(ave_grads)), ave_grads, align="center"
10            )
11     plt.xticks(range(len(ave_grads)), layers, rotation=90)
12     plt.xlabel("Layers")
13     plt.ylabel("Average gradient")
14     plt.title("Gradient flow")
15     plt.tight_layout()
16     plt.show()
17
18 # Usage
19 plot_grad_flow(model.named_parameters())
```

Summary

- **Automatic differentiation** provides efficient gradient computation for neural network training
- **Computational graphs** form the conceptual and practical framework for implementing autograd
- **Vanishing gradients** represent a fundamental challenge in training deep networks
- Modern solutions include:
 - Better activation functions (ReLU family)
 - Improved weight initialization
 - Architectural innovations (skip connections, normalization)
- Understanding these concepts is crucial for designing and debugging deep neural networks

- Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. CVPR.
- Baydin, A. G., Pearlmutter, B. A., Radul, A. A., & Siskind, J. M. (2018). Automatic differentiation in machine learning: a survey. JMLR.
- Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. ICML.
- PyTorch documentation:
<https://pytorch.org/docs/stable/autograd.html>