

# Summary of Activation Functions, Loss Functions, and Optimizers

Md Shahriar Forhad (<https://github.com/Shahriar88>)

# Contents

<b>1</b>	<b>Activation Functions, Loss Functions, Optimizers, and Training Mechanics</b>	<b>1</b>
1.1	Activation Functions	1
1.1.1	Sigmoid	1
1.1.2	Tanh	1
1.1.3	ReLU	2
1.1.4	Leaky ReLU	2
1.1.5	Softmax	2
1.1.6	GELU	2
1.2	Loss Functions	2
1.2.1	Mean Squared Error (MSE)	2
1.2.2	Mean Absolute Error (MAE)	2
1.2.3	Huber Loss	3
1.2.4	Cross Entropy Loss	3
1.2.5	Focal Loss	3
1.2.6	Dice / IoU Loss	3
1.2.7	KL Divergence	3
1.3	Optimizers	4
1.3.1	Stochastic Gradient Descent (SGD)	4
1.3.2	Momentum	5
1.3.3	Adagrad	5
1.3.4	RMSProp	5
1.3.5	Adam	5
1.3.6	AdamW	5
1.3.7	LAMB	5
1.4	Summary Tables	6
1.5	Forward Pass	7
1.6	Backpropagation	7
1.7	Weight and Bias Updates	7
<b>2</b>	<b>General Structure of Forward Pass and Backpropagation</b>	<b>9</b>
2.1	Overview	9
2.2	Forward Pass	9
2.2.1	Basic Structure	9
2.2.2	Computation Flow	9
2.3	Backpropagation	9
2.3.1	Error Propagation	9
2.3.2	Weight Update	10
2.3.3	General Steps	10
2.4	Universality of the Process	10
2.5	Common Variations	10
2.6	Summary	10

2.7	Worked Example I: Single-Neuron Forward and Backpropagation . . . . .	11
2.8	Worked Example II: Three-Layer Forward and Backpropagation . . . . .	11
2.8.1	Network Definition (Scalar) . . . . .	11
2.8.2	Activation Derivatives . . . . .	12
2.8.3	Numerical Values Used . . . . .	12
2.9	Matrix-Form Backpropagation, Product Types, and Forward-Pass Memory . . .	15
2.10	Chain Rule Accumulation and Gradient Flow . . . . .	17
2.11	Hessian–Vector Products and Second-Order Structure . . . . .	17
2.11.1	Worked Example: Scalar Hessian–Vector Product . . . . .	18
2.11.2	Hessian–Vector Products via Backpropagation: Concrete Three-Layer Ex- ample . . . . .	18
2.11.3	Role of the Second Backward Pass . . . . .	24
2.11.4	Direction Vectors in Second-Order Optimization . . . . .	24
2.11.5	Newton–CG and Hessian-Free Methods . . . . .	24
2.12	Training vs. Inference in Automatic Differentiation Frameworks . . . . .	26
<b>3</b>	<b>Jacobians in Backpropagation and Second-Order Analysis</b>	<b>27</b>
3.1	Overview . . . . .	27
3.2	Jacobian View of the Chain Rule . . . . .	27
3.3	Why Jacobians Are Not Formed Explicitly . . . . .	27
3.4	Three-Layer Network: Matrix Form . . . . .	28
3.5	Numerical Example with Explicit Jacobians (with activation and loss derivatives)	28
3.5.1	Network definition . . . . .	28
3.5.2	Activation derivatives and loss derivative . . . . .	29
3.5.3	Forward pass (numerical) . . . . .	29
3.5.4	Local derivatives (numerical) . . . . .	29
3.5.5	Explicit Jacobians . . . . .	30
3.5.6	Backpropagation via Jacobian products (numerical) . . . . .	30
3.6	Connection to Second-Order Derivatives . . . . .	31
3.7	Summary . . . . .	31
<b>4</b>	<b>Vision Transformer (ViT)</b>	<b>33</b>
4.1	Overview . . . . .	33
4.2	Patch Embedding and Tokenization . . . . .	33
4.3	Transformer Depth and Attention Modules . . . . .	33
4.4	Multi-Head Self-Attention . . . . .	34
4.5	Computational Scaling . . . . .	34
4.5.1	Attention Cost . . . . .	34
4.5.2	MLP Cost . . . . .	34
4.6	Parameter Impact Summary . . . . .	34
4.7	Numerical Example . . . . .	35
4.8	Implications for Detection and Segmentation . . . . .	35
4.9	Summary . . . . .	35

# List of Figures

4.1	High-level architecture of the Vision Transformer (ViT). An input image is split into non-overlapping patches, embedded into a token sequence, and processed by a Transformer encoder stack of depth $L$ . Each encoder layer contains exactly one multi-head self-attention (MHSA) module (so the number of attention modules is $L$ ), and each MHSA uses $h$ attention heads. . . . .	36
4.2	Hydra/tree-style view of ViT without overflow. Tokens pass through an encoder layer repeated $\times L$ . Each layer contains exactly one MHSA (so the number of attention modules is $L$ ) with $h$ heads and one MLP block. The MHSA splits into $h$ heads and merges back via concatenation to dimension $D$ . . . . .	37
4.3	Tree-style view of ViT. The encoder stack has depth $L$ ; each layer contains one MHSA (therefore total attention modules = $L$ ) with $h$ heads and an MLP block. . . . .	38



# List of Tables

1.1	Forward Pass Computations and Data Structures . . . . .	7
1.2	Comprehensive Backpropagation Steps . . . . .	8
1.3	Weight and Bias Update Steps . . . . .	8
2.1	Forward propagation for a single-neuron network . . . . .	11
2.2	Backpropagation for the output layer (layer 2) . . . . .	11
2.3	Backpropagation for the hidden layer (layer 1) . . . . .	12
2.4	Forward propagation for a three-layer network with three different activation functions . . . . .	12
2.5	Backpropagation through the output layer (layer 3, sigmoid) . . . . .	13
2.6	Backpropagation through hidden layer 2 (layer 2, tanh) . . . . .	13
2.7	Backpropagation through hidden layer 1 (layer 1, ReLU) . . . . .	13
2.8	Quantities stored during forward propagation and their role in backpropagation .	16
2.9	Product types appearing in forward and backward propagation . . . . .	16
2.10	Comparison of PyTorch behavior during training and inference . . . . .	25
4.1	Impact of Vision Transformer architectural parameters. . . . .	35
4.2	Impact of the MLP ratio in Vision Transformers. . . . .	36



# Chapter 1

## Activation Functions, Loss Functions, Optimizers, and Training Mechanics

### 1.1 Activation Functions

#### 1.1.1 Sigmoid

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

**Applications:** Binary classification outputs, logistic regression, simple neural networks. **Advantages:**

- Smooth, differentiable.
- Maps input to  $(0, 1)$  range for probability interpretation.

**Disadvantages:**

- Vanishing gradients for large  $|x|$ .
- Non-zero-centered outputs slow convergence.

#### 1.1.2 Tanh

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

**Applications:** Hidden layers of RNNs, MLPs. **Advantages:**

- Zero-centered outputs.
- Stronger gradients than Sigmoid.

**Disadvantages:**

- Still suffers from vanishing gradients.
- Costlier than ReLU.

### 1.1.3 ReLU

$$\text{ReLU}(x) = \max(0, x)$$

**Applications:** CNNs, fully connected layers in deep networks. **Advantages:**

- Simple and efficient.
- Sparse activations reduce computation.
- Mitigates vanishing gradient problem.

**Disadvantages:**

- “Dying ReLU” problem for negative inputs.
- Unbounded positive outputs.

### 1.1.4 Leaky ReLU

$$f(x) = \begin{cases} x, & x > 0 \\ \alpha x, & x \leq 0 \end{cases}$$

**Applications:** CNNs where ReLU causes dead neurons. **Advantages:** Avoids dying ReLU; small slope for negatives. **Disadvantages:** Requires  $\alpha$  tuning.

### 1.1.5 Softmax

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

**Applications:** Multiclass classification output layers. **Advantages:** Converts logits to probabilities. **Disadvantages:** Sensitive to large logits; saturation slows learning.

### 1.1.6 GELU

$$\text{GELU}(x) = x\Phi(x)$$

**Applications:** Transformers, BERT, Vision Transformers. **Advantages:** Smooth and probabilistic activation. **Disadvantages:** More computation-heavy.

## 1.2 Loss Functions

### 1.2.1 Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

**Applications:** Regression problems, autoencoders. **Advantages:** Penalizes large errors heavily; convex. **Disadvantages:** Sensitive to outliers.

### 1.2.2 Mean Absolute Error (MAE)

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

**Applications:** Robust regression. **Advantages:** Robust to outliers. **Disadvantages:** Non-differentiable at 0.

### 1.2.3 Huber Loss

$$L_{\delta}(e) = \begin{cases} \frac{1}{2}e^2, & |e| \leq \delta \\ \delta(|e| - \frac{1}{2}\delta), & |e| > \delta \end{cases}$$

**Applications:** Regression with moderate outliers. **Advantages:** Combines smoothness (MSE) and robustness (MAE). **Disadvantages:** Requires tuning  $\delta$ .

### 1.2.4 Cross Entropy Loss

$$L = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log p_{ik}$$

**Applications:** Classification (binary, multiclass, multi-label). **Advantages:** Probabilistic; aligns with maximum likelihood. **Disadvantages:** Sensitive to noisy labels; overconfident predictions.

### 1.2.5 Focal Loss

$$L = -\alpha(1 - p_t)^{\gamma} \log(p_t)$$

**Applications:** Object detection (Mask R-CNN, RetinaNet). **Advantages:** Handles class imbalance. **Disadvantages:** Hyperparameters  $\alpha, \gamma$  need tuning.

### 1.2.6 Dice / IoU Loss

$$\text{Dice} = \frac{2|A \cap B|}{|A| + |B|}, \quad L = 1 - \text{Dice}$$

**Applications:** Image segmentation (Mask R-CNN, U-Net). **Advantages:** Works well with imbalanced masks. **Disadvantages:** Non-linear, harder optimization.

### 1.2.7 KL Divergence

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

**Applications:** Variational Autoencoders, distillation. **Advantages:** Measures information difference. **Disadvantages:** Asymmetric; can be unstable.

## Loss Functions

Loss	Typical Use	Advantages	Disadvantages
MSE ( $\frac{1}{N} \sum (y - \hat{y})^2$ )	Regression	Convex, smooth, strong penalty on large errors	Sensitive to outliers
MAE ( $\frac{1}{N} \sum  y - \hat{y} $ )	Robust regression	Robust to outliers	Non-differentiable at 0; slower convergence
Huber / Smooth- $L_1$	Regression, box regression	Combines MSE and MAE; robust and smooth	Requires tuning of $\delta$
Binary CE	Binary classification	Probabilistic; aligns with MLE	Sensitive to noisy labels
Multiclass CE	Multiclass classification	Standard loss with softmax; probabilistic interpretation	Overconfident predictions are heavily penalized
BCEWithLogits	Multi-label classification	Numerically stable (sigmoid + CE combined)	Ignores label dependencies
Focal Loss	Imbalanced classification/detection	Focuses learning on hard examples	Requires tuning of $\alpha$ and $\gamma$
Dice Loss	Image segmentation (masks)	Handles class imbalance; measures overlap directly	Unstable for very small targets
IoU (Jaccard) Loss	Segmentation / bounding boxes	Directly optimizes IoU metric	Non-smooth; slower early convergence
Smooth- $L_1$ (boxes)	Object detection (R-CNNs)	Robust to outliers; standard for box regression	Scale-sensitive
KL Divergence	VAEs, knowledge distillation	Measures divergence between distributions	Asymmetric; unstable when $Q \approx 0$
Triplet Loss	Metric learning	Learns discriminative embedding spaces	Needs triplet mining; margin tuning
Contrastive Loss	Siamese networks / similarity tasks	Learns pairwise distance relationships	Requires balanced positive/negative pairs
Cosine Embedding	Text/vision embeddings	Scale-invariant and simple	Loses magnitude information
Perceptual Loss	Super-resolution / style transfer	Encourages perceptual similarity using deep features	Computationally heavy; needs pretrained $\phi$
Total Variation	Image smoothing / denoising	Removes noise, encourages smoothness	Can over-smooth fine details

## 1.3 Optimizers

### 1.3.1 Stochastic Gradient Descent (SGD)

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} J(\theta_t)$$

**Applications:** Classical ML, CNNs. **Advantages:** Simple, effective. **Disadvantages:** Sensitive to learning rate; oscillates.

### 1.3.2 Momentum

$$v_t = \beta v_{t-1} + (1 - \beta) \nabla_{\theta} J(\theta_t), \quad \theta_{t+1} = \theta_t - \eta v_t$$

**Applications:** Deep CNNs. **Advantages:** Faster convergence; less noise. **Disadvantages:** Needs tuning of  $\beta$ .

### 1.3.3 Adagrad

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot \nabla_{\theta} J(\theta_t)$$

**Applications:** Sparse features (e.g. NLP embeddings). **Advantages:** Adaptive learning rate. **Disadvantages:** Learning rate decays too fast.

### 1.3.4 RMSProp

$$E[g^2]_t = \beta E[g^2]_{t-1} + (1 - \beta) g_t^2, \quad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

**Applications:** RNNs, time series, non-stationary data. **Advantages:** Stable, adaptive. **Disadvantages:** Sensitive to  $\beta$ .

### 1.3.5 Adam

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t, \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}, \quad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$

**Applications:** Deep learning, NLP, Transformers. **Advantages:** Combines momentum + adaptive rate. **Disadvantages:** May generalize poorly.

### 1.3.6 AdamW

$$\theta_{t+1} = \theta_t - \eta \left( \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} + \lambda \theta_t \right)$$

**Applications:** Transformers, BERT, ViTs. **Advantages:** Decoupled weight decay  $\rightarrow$  better generalization. **Disadvantages:** Slightly more computation.

### 1.3.7 LAMB

$$r_t = \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}, \quad \theta_{t+1} = \theta_t - \eta \frac{\|\theta_t\|}{\|r_t\|} r_t$$

**Applications:** Large-batch Transformer training. **Advantages:** Enables distributed training. **Disadvantages:** Complex; more tuning.

## 1.4 Summary Tables

### Activation Functions

Function	Pros	Cons
Sigmoid	Probabilistic, smooth	Vanishing gradient
Tanh	Zero-centered	Saturation at extremes
ReLU	Sparse, fast	Dead neurons
Leaky ReLU	Fixes dead ReLU	Hyperparam $\alpha$
Softmax	Probabilities	Saturation
GELU	Smooth	Slower compute

### Optimizers

Optimizer	Pros	Cons
SGD	Simple, reliable	Slow, oscillates
Momentum	Smooth convergence	Needs tuning
RMSProp	Stable	Sensitive $\beta$
Adam	Fast, adaptive	May overfit
AdamW	Best generalization	Slightly slower

## 1.5 Forward Pass

The forward pass computes activations layer by layer using learned weights and biases. Each layer applies linear transformations followed by nonlinear activation functions (e.g., sigmoid, ReLU, etc.).

Table 1.1: Forward Pass Computations and Data Structures

Step	Calculation	Data Structure	Multiplication Type
<b>Forward Pass</b>			
Input to 1st hidden layer (HI_1)	$HI_1 = W_1 \cdot X + B_1$	$HI_1$ : Matrix, $W_1$ : Matrix, $X$ : Matrix, $B_1$ : Vector	<b>Dot Product</b> (between $W_1$ and $X$ )
Output of 1st hidden layer (HO_1)	$HO_1 = \sigma(HI_1)$	$HO_1$ : Vector	<b>Element-wise</b> (Sigmoid applied element-wise)
Input to 2nd hidden layer (HI_2)	$HI_2 = W_2 \cdot HO_1 + B_2$	$HI_2$ : Vector, $W_2$ : Matrix, $HO_1$ : Vector, $B_2$ : Vector	<b>Dot Product</b> (between $W_2$ and $HO_1$ )
Output of 2nd hidden layer (HO_2)	$HO_2 = \sigma(HI_2)$	$HO_2$ : Vector	<b>Element-wise</b> (Sigmoid applied element-wise)
Input to output layer (HO_final)	$HO_{final} = W_3 \cdot HO_2 + B_3$	$HO_{final}$ : Vector, $W_3$ : Matrix, $HO_2$ : Vector, $B_3$ : Vector	<b>Dot Product</b> (between $W_3$ and $HO_2$ )
Final output $\hat{Y}$	$\hat{Y} = \sigma(HO_{final})$	$\hat{Y}$ : Vector	<b>Element-wise</b> (Sigmoid applied element-wise)
Error Calculation	$E = Y - \hat{Y}$	$E$ : Vector, $Y$ : Vector, $\hat{Y}$ : Vector	<b>Element-wise</b> (Subtraction)

## 1.6 Backpropagation

$$Error\_Layer\_N = \sigma'_d(Input\_N) \left( Error\_Layer\_N + 1 (W\_N + 1)^T \right), \quad Input\_N = WX + B.$$

## 1.7 Weight and Bias Updates

$$W_{new} = W_{old} + lr \cdot Own\_Error \cdot Own\_Input^T$$

$$Own\_Input = X, \quad Own\_Input\_N = HO_{N-1} = \sigma(HI_{N-1})$$

Table 1.2: Comprehensive Backpropagation Steps

Step	Calculation	Data Structure	Multiplication Type
<b>Backpropagation</b>			
Error at output layer	$err_{HO_{final}} = E \cdot \sigma'(HO_{final})$	$err_{HO_{final}}$ : Vector, $E$ : Vector, $\sigma'(HO_{final})$ : Vector	<b>Element-wise</b> (Multiplication)
Error in 2nd hidden layer (err_HO_2)	$err_{HO_2} = err_{HO_{final}} \cdot W_3^T \cdot \sigma'(HI_2)$	$err_{HO_2}$ : Vector, $W_3^T$ : Matrix (transpose), $\sigma'(HI_2)$ : Vector	<b>Dot Product</b> (between $err_{HO_{final}}$ and $W_3^T$ ) followed by <b>Element-wise</b> multiplication with $\sigma'(HI_2)$
Error in 1st hidden layer (err_HO_1)	$err_{HO_1} = err_{HO_2} \cdot W_2^T \cdot \sigma'(HI_1)$	$err_{HO_1}$ : Vector, $W_2^T$ : Matrix (transpose), $\sigma'(HI_1)$ : Vector	<b>Dot Product</b> (between $err_{HO_2}$ and $W_2^T$ ) followed by <b>Element-wise</b> multiplication with $\sigma'(HI_1)$
Error Calculation	$E = Y - \hat{Y}$	$E$ : Vector, $Y$ : Vector, $\hat{Y}$ : Vector	<b>Element-wise</b> (Subtraction)

Table 1.3: Weight and Bias Update Steps

Step	Calculation	Data Structure	Multiplication Type
<b>Weight and Bias Updates</b>			
Update $W_1$ and $B_1$	$W_1 = W_1 + lr \cdot X^T \cdot err_{HO_1}$ , $B_1 = B_1 + lr \cdot err_{HO_1}$	$W_1$ : Matrix, $X^T$ : Matrix (transpose), $err_{HO_1}$ : Vector	<b>Dot Product</b> (between $X^T$ and $err_{HO_1}$ )
Update $W_2$ and $B_2$	$W_2 = W_2 + lr \cdot HO_1^T \cdot err_{HO_2}$ , $B_2 = B_2 + lr \cdot err_{HO_2}$	$W_2$ : Matrix, $HO_1^T$ : Matrix (transpose), $err_{HO_2}$ : Vector	<b>Dot Product</b> (between $HO_1^T$ and $err_{HO_2}$ )
Update $W_3$ and $B_3$	$W_3 = W_3 + lr \cdot HO_2^T \cdot err_{HO_{final}}$ , $B_3 = B_3 + lr \cdot err_{HO_{final}}$	$W_3$ : Matrix, $HO_2^T$ : Matrix (transpose), $err_{HO_{final}}$ : Vector	<b>Dot Product</b> (between $HO_2^T$ and $err_{HO_{final}}$ )

## Chapter 2

# General Structure of Forward Pass and Backpropagation

### 2.1 Overview

The general structure of the forward pass and backpropagation is the same for all basic neural networks, regardless of how many hidden layers or neurons are used. The following principles hold true for most feedforward neural networks (also known as **multilayer perceptrons (MLPs)**).

### 2.2 Forward Pass

#### 2.2.1 Basic Structure

- The network consists of an input layer, one or more hidden layers, and an output layer.
- Each layer computes a weighted sum of its inputs, adds a bias, and applies an activation function (such as sigmoid, ReLU, or tanh) to produce its output.

#### 2.2.2 Computation Flow

- Data flows forward through the network from the input layer to the output layer.
- The output of one layer becomes the input to the next layer.

### 2.3 Backpropagation

#### 2.3.1 Error Propagation

- After computing the final output during the forward pass, the network compares the predicted output with the expected (target) output.
- The error is calculated using a loss function (for example, mean squared error or cross-entropy).
- This error is propagated backward from the output layer to the hidden layers, adjusting the neuron weights to minimize the error.

### 2.3.2 Weight Update

- The gradients (rate of change of error with respect to weights) are computed using the chain rule of calculus.
- The weights and biases are updated using gradient descent or its variants (e.g., stochastic gradient descent).

### 2.3.3 General Steps

1. Compute the error at the output layer.
2. Backpropagate the error through each preceding layer.
3. Update the weights and biases to reduce the overall error.

## 2.4 Universality of the Process

The process of forward pass and backpropagation described above holds for all basic feedforward neural networks, regardless of:

- The number of hidden layers.
- The number of neurons per layer.
- The activation function used.

## 2.5 Common Variations

While the core steps remain the same, there are variations among different network types and techniques:

- **Activation Functions:** Modern networks often use ReLU or tanh instead of sigmoid in hidden layers.
- **Loss Functions:** Classification problems typically use cross-entropy loss, while regression problems often use mean squared error (MSE).
- **Learning Algorithms:** Standard gradient descent can be replaced by advanced optimizers such as Adam, RMSProp, or Adagrad.

## 2.6 Summary

- The concept of forward pass and backpropagation is **universal** to most feedforward neural networks.
- Differences arise mainly in the architecture (number of layers or neurons) and the type of activation, loss function, and optimization algorithm used.
- As networks become deeper and more complex, the same principles apply, but across more layers.

### Examples

- **Deep Neural Networks (DNNs):** Contain more hidden layers, but the core forward and backward propagation remain identical.
- **Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs):** Both employ forward and backward propagation with specialized modifications.

In essence, for any basic feedforward neural network, the process of forward pass and backpropagation remains fundamentally the same.

## 2.7 Worked Example I: Single-Neuron Forward and Backpropagation

Table 2.1: Forward propagation for a single-neuron network

Step	Quantity	Expression	Numerical value
1	Input	$x$	2
2	Weight	$w$	3
3	Bias	$b$	1
4	Linear output	$z = wx + b$	7
5	Activation	$\hat{y} = \sigma(z)$	0.999088
6	Target	$y$	1
7	Loss	$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2$	$4.16 \times 10^{-7}$

Table 2.2: Backpropagation for the output layer (layer 2)

Step	Gradient	Formula	Value	Value comes from
7	$\frac{\partial \mathcal{L}}{\partial \hat{y}}$	$\hat{y} - y$	=0.119587	Forward step 5 ( $\hat{y}$ ), forward step 6 ( $y$ )
8	$\frac{\partial \hat{y}}{\partial z_2}$	$\hat{y}(1 - \hat{y})$	0.105233	Forward step 5 ( $\hat{y}$ )
9	$\frac{\partial \mathcal{L}}{\partial z_2}$	Step 7 $\times$ Step 8	=0.012580	Backprop steps 7 and 8
10	$\frac{\partial z_2}{\partial w_2}$	$a_1$	0.999088	Forward step 3 ( $a_1$ )
11	$\frac{\partial \mathcal{L}}{\partial w_2}$	Step 9 $\times$ Step 10	=0.012569	Backprop step 9, forward step 3
12	$\frac{\partial z_2}{\partial b_2}$	1	1	Affine layer definition
13	$\frac{\partial \mathcal{L}}{\partial b_2}$	Step 9 $\times$ Step 12	=0.012580	Backprop step 9

**Note.** Each gradient is computed using the chain rule and depends exclusively on quantities obtained during the forward pass or gradients propagated from subsequent layers. This explicit dependency structure mirrors the computation graph used by automatic differentiation frameworks.

## 2.8 Worked Example II: Three-Layer Forward and Backpropagation

### 2.8.1 Network Definition (Scalar)

$$z_1 = w_1x + b_1, \quad a_1 = \text{ReLU}(z_1),$$

Table 2.3: Backpropagation for the hidden layer (layer 1)

Step	Gradient	Formula	Value	Value comes from
14	$\frac{\partial z_2}{\partial a_1}$	$w_2$	4	Forward definition of layer 2
15	$\frac{\partial \mathcal{L}}{\partial a_1}$	Step 9 $\times$ Step 14	=0.050319	Backprop step 9, forward step 4
16	$\frac{\partial a_1}{\partial z_1}$	$a_1(1 - a_1)$	0.000911	Forward step 3 ( $a_1$ )
17	$\frac{\partial \mathcal{L}}{\partial z_1}$	Step 15 $\times$ Step 16	= $4.584 \times 10^{-5}$	Backprop steps 15 and 16
18	$\frac{\partial z_1}{\partial w_1}$	$x$	2	Forward step 1 ( $x$ )
19	$\frac{\partial \mathcal{L}}{\partial w_1}$	Step 17 $\times$ Step 18	= $9.168 \times 10^{-5}$	Backprop step 17, forward step 1
20	$\frac{\partial z_1}{\partial b_1}$	1	1	Affine layer definition
21	$\frac{\partial \mathcal{L}}{\partial b_1}$	Step 17 $\times$ Step 20	= $4.584 \times 10^{-5}$	Backprop step 17

$$z_2 = w_2 a_1 + b_2, \quad a_2 = \tanh(z_2),$$

$$z_3 = w_3 a_2 + b_3, \quad \hat{y} = \sigma(z_3), \quad \mathcal{L} = \frac{1}{2}(y - \hat{y})^2.$$

### 2.8.2 Activation Derivatives

$$\frac{d}{dz} \text{ReLU}(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}, \quad \frac{d}{dz} \tanh(z) = 1 - \tanh^2(z), \quad \frac{d}{dz} \sigma(z) = \sigma(z)(1 - \sigma(z)).$$

### 2.8.3 Numerical Values Used

$$x = 2, \quad (w_1, b_1) = (1.5, -1), \quad (w_2, b_2) = (0.5, 0.1), \quad (w_3, b_3) = (2, -0.3), \quad y = 1.$$

Table 2.4: Forward propagation for a three-layer network with three different activation functions

Step	Quantity	Expression	Numerical value
1	Input	$x$	2
2	Layer 1 weight	$w_1$	1.5
3	Layer 1 bias	$b_1$	=1
4	Pre-activation (L1)	$z_1 = w_1 x + b_1$	2.0
5	Activation (L1, ReLU)	$a_1 = \text{ReLU}(z_1)$	2.0
6	Layer 2 weight	$w_2$	0.5
7	Layer 2 bias	$b_2$	0.1
8	Pre-activation (L2)	$z_2 = w_2 a_1 + b_2$	1.1
9	Activation (L2, tanh)	$a_2 = \tanh(z_2)$	0.800499
10	Layer 3 weight	$w_3$	2
11	Layer 3 bias	$b_3$	=0.3
12	Pre-activation (L3)	$z_3 = w_3 a_2 + b_3$	1.300998
13	Output activation (sigmoid)	$\hat{y} = \sigma(z_3)$	0.786003
14	Target	$y$	1
15	Loss	$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2$	0.022897

Table 2.5: Backpropagation through the output layer (layer 3, sigmoid)

Step	Gradient	Formula	Value	Value comes from
16	$\frac{\partial \mathcal{L}}{\partial \hat{y}}$	$\hat{y} - y$	=0.213997	Forward step 13 ( $\hat{y}$ ), step 14 ( $y$ )
17	$\frac{\partial \hat{y}}{\partial z_3}$	$\hat{y}(1 - \hat{y})$	0.168202	Forward step 13 ( $\hat{y}$ )
18	$\frac{\partial \mathcal{L}}{\partial z_3}$	Step 16 $\times$ Step 17	=0.035995	Backprop steps 16 and 17
19	$\frac{\partial z_3}{\partial w_3}$	$a_2$	0.800499	Forward step 9 ( $a_2$ )
20	$\frac{\partial \mathcal{L}}{\partial w_3}$	Step 18 $\times$ Step 19	=0.028814	Backprop step 18, forward step 9
21	$\frac{\partial z_3}{\partial b_3}$	1	1	Affine layer definition
22	$\frac{\partial \mathcal{L}}{\partial b_3}$	Step 18 $\times$ Step 21	=0.035995	Backprop step 18
23	$\frac{\partial z_3}{\partial a_2}$	$w_3$	2	Forward step 10 ( $w_3$ )
24	$\frac{\partial \mathcal{L}}{\partial a_2}$	Step 18 $\times$ Step 23	=0.071990	Backprop step 18, forward step 10

Table 2.6: Backpropagation through hidden layer 2 (layer 2, tanh)

Step	Gradient	Formula	Value	Value comes from
25	$\frac{\partial a_2}{\partial z_2}$	$1 - a_2^2$	0.359201	Forward step 9 ( $a_2$ )
26	$\frac{\partial \mathcal{L}}{\partial z_2}$	Step 24 $\times$ Step 25	=0.025859	Backprop step 24, backprop step 25
27	$\frac{\partial z_2}{\partial w_2}$	$a_1$	2.0	Forward step 5 ( $a_1$ )
28	$\frac{\partial \mathcal{L}}{\partial w_2}$	Step 26 $\times$ Step 27	=0.051718	Backprop step 26, forward step 5
29	$\frac{\partial z_2}{\partial b_2}$	1	1	Affine layer definition
30	$\frac{\partial \mathcal{L}}{\partial b_2}$	Step 26 $\times$ Step 29	=0.025859	Backprop step 26
31	$\frac{\partial z_2}{\partial a_1}$	$w_2$	0.5	Forward step 6 ( $w_2$ )
32	$\frac{\partial \mathcal{L}}{\partial a_1}$	Step 26 $\times$ Step 31	=0.012929	Backprop step 26, forward step 6

Table 2.7: Backpropagation through hidden layer 1 (layer 1, ReLU)

Step	Gradient	Formula	Value	Value comes from
33	$\frac{\partial a_1}{\partial z_1}$	$\mathbb{I}[z_1 > 0]$	1	Forward step 4 ( $z_1 = 2.0 > 0$ )
34	$\frac{\partial \mathcal{L}}{\partial z_1}$	Step 32 $\times$ Step 33	=0.012929	Backprop step 32, backprop step 33
35	$\frac{\partial z_1}{\partial w_1}$	$x$	2	Forward step 1 ( $x$ )
36	$\frac{\partial \mathcal{L}}{\partial w_1}$	Step 34 $\times$ Step 35	=0.025859	Backprop step 34, forward step 1
37	$\frac{\partial z_1}{\partial b_1}$	1	1	Affine layer definition
38	$\frac{\partial \mathcal{L}}{\partial b_1}$	Step 34 $\times$ Step 37	=0.012929	Backprop step 34

**Note.** Each gradient is computed using the chain rule and depends exclusively on quantities obtained during the forward pass or gradients propagated from subsequent layers. The three activation derivatives appear explicitly: sigmoid uses  $\hat{y}(1 - \hat{y})$ , tanh uses  $1 - a_2^2$ , and ReLU uses the gate  $\mathbb{K}[z_1 > 0]$ .

## 2.9 Matrix-Form Backpropagation, Product Types, and Forward-Pass Memory

**Goal.** This section presents forward and backward propagation for a three-layer feedforward neural network using vector and matrix notation. The formulation explicitly identifies the type of each algebraic operation (matrix–vector product, Hadamard product, outer product) and clarifies which intermediate quantities must be stored during the forward pass to enable backpropagation.

**Notation and dimensions.** Let the input be  $\mathbf{x} \in \mathbb{R}^{d_0}$  and define  $\mathbf{a}_0 := \mathbf{x}$ . For layers  $k \in \{1, 2, 3\}$ :

$$\mathbf{W}_k \in \mathbb{R}^{d_k \times d_{k-1}}, \quad \mathbf{b}_k \in \mathbb{R}^{d_k}, \quad \mathbf{z}_k \in \mathbb{R}^{d_k}, \quad \mathbf{a}_k \in \mathbb{R}^{d_k}.$$

Each affine transformation uses a **matrix–vector product**:

$$\mathbf{z}_k = \mathbf{W}_k \mathbf{a}_{k-1} + \mathbf{b}_k.$$

**Three-layer network with heterogeneous activations.**

$$\mathbf{a}_1 = \text{ReLU}(\mathbf{z}_1), \quad \mathbf{a}_2 = \tanh(\mathbf{z}_2), \quad \hat{\mathbf{y}} = \mathbf{a}_3 = \sigma(\mathbf{z}_3).$$

For a single training sample, the loss is

$$\mathcal{L} = \frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{y}\|^2.$$

**Activation derivatives (element-wise).**

$$\frac{d}{dz} \text{ReLU}(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}, \quad \frac{d}{dz} \tanh(z) = 1 - \tanh^2(z), \quad \frac{d}{dz} \sigma(z) = \sigma(z)(1 - \sigma(z)).$$

The symbol  $\odot$  denotes the *Hadamard (element-wise) product*.

**Forward propagation (matrix form).**

$$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1, \quad \mathbf{a}_1 = \text{ReLU}(\mathbf{z}_1), \quad (2.1)$$

$$\mathbf{z}_2 = \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2, \quad \mathbf{a}_2 = \tanh(\mathbf{z}_2), \quad (2.2)$$

$$\mathbf{z}_3 = \mathbf{W}_3 \mathbf{a}_2 + \mathbf{b}_3, \quad \hat{\mathbf{y}} = \sigma(\mathbf{z}_3). \quad (2.3)$$

**Backpropagation (matrix form).** Define error signals (“deltas”):

$$\boldsymbol{\delta}_k := \frac{\partial \mathcal{L}}{\partial \mathbf{z}_k} \in \mathbb{R}^{d_k}.$$

**Layer 3 (sigmoid output).**

$$\boldsymbol{\delta}_3 = (\hat{\mathbf{y}} - \mathbf{y}) \odot (\hat{\mathbf{y}} \odot (1 - \hat{\mathbf{y}})).$$

**Layer 2 (tanh).**

$$\boldsymbol{\delta}_2 = (\mathbf{W}_3^\top \boldsymbol{\delta}_3) \odot (1 - \mathbf{a}_2^{\odot 2}).$$

**Layer 1 (ReLU).**

$$\boldsymbol{\delta}_1 = (\mathbf{W}_2^\top \boldsymbol{\delta}_2) \odot \mathbb{I}[\mathbf{z}_1 > 0].$$

Table 2.8: Quantities stored during forward propagation and their role in backpropagation

Stored quantity	Saved during forward pass	Used in backpropagation
$\mathbf{x} = \mathbf{a}_0$	Input layer	$\boldsymbol{\delta}_1 \mathbf{x}^\top$
$\mathbf{z}_1$	Pre-activation (layer 1)	ReLU mask $\mathbb{K}[\mathbf{z}_1 > 0]$
$\mathbf{a}_1$	Activation (layer 1)	$\boldsymbol{\delta}_2 \mathbf{a}_1^\top$
$\mathbf{z}_2$	Pre-activation (layer 2)	$1 - \mathbf{a}_2^{\odot 2}$
$\mathbf{a}_2$	Activation (layer 2)	$\boldsymbol{\delta}_3 \mathbf{a}_2^\top$
$\mathbf{z}_3$	Pre-activation (output layer)	$\hat{\mathbf{y}} \odot (1 - \hat{\mathbf{y}})$
$\hat{\mathbf{y}}$	Network output	$\hat{\mathbf{y}} - \mathbf{y}$

**Parameter gradients.** For each layer  $k \in \{1, 2, 3\}$ :

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_k} = \boldsymbol{\delta}_k \mathbf{a}_{k-1}^\top, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}_k} = \boldsymbol{\delta}_k.$$

The product  $\boldsymbol{\delta}_k \mathbf{a}_{k-1}^\top$  is an **outer product**.

**Forward-pass memory (saved tensors).**

Table 2.9: Product types appearing in forward and backward propagation

Symbol / form	Name	Where it appears
$\mathbf{W}\mathbf{a}, \mathbf{W}^\top \boldsymbol{\delta}$	Matrix–vector product	Forward affine map; backward error propagation
$\odot$	Hadamard (element-wise) product	Activation derivatives and masks
$\boldsymbol{\delta} \mathbf{a}^\top$	Outer product	Weight gradients

**Summary of product types.**

## 2.10 Chain Rule Accumulation and Gradient Flow

**Motivation.** In backpropagation, gradients are accumulated whenever a variable influences the loss through multiple computational paths. This section makes the summation in the chain rule explicit and connects it directly to how automatic differentiation frameworks accumulate gradients.

**Computation graph example.** Consider the scalar computation

$$\begin{aligned} v_3 &= x^2, \\ v_5 &= 2v_3, \\ v_6 &= v_3 + 1, \\ f &= v_5 v_6. \end{aligned}$$

The intermediate variable  $v_3$  feeds into two downstream nodes.

**Forward pass.** For  $x = 3$ :

$$v_3 = 9, \quad v_5 = 18, \quad v_6 = 10, \quad f = 180.$$

**Backward pass.** The multivariable chain rule gives

$$\frac{\partial f}{\partial v_3} = \frac{\partial f}{\partial v_6} \frac{\partial v_6}{\partial v_3} + \frac{\partial f}{\partial v_5} \frac{\partial v_5}{\partial v_3}.$$

Evaluating:

$$\begin{aligned} \frac{\partial f}{\partial v_5} &= v_6 = 10, & \frac{\partial f}{\partial v_6} &= v_5 = 18, \\ \frac{\partial v_5}{\partial v_3} &= 2, & \frac{\partial v_6}{\partial v_3} &= 1. \end{aligned}$$

Thus,

$$\frac{\partial f}{\partial v_3} = 18 + 20 = 38.$$

**Interpretation.** The summation corresponds exactly to gradient accumulation in backpropagation: each outgoing edge contributes a partial gradient, and the total gradient is their sum.

## 2.11 Hessian–Vector Products and Second-Order Structure

Let  $f(\boldsymbol{\theta})$  be a scalar loss with parameters  $\boldsymbol{\theta} \in \mathbb{R}^W$ . The gradient is  $\nabla f(\boldsymbol{\theta}) \in \mathbb{R}^W$  and the Hessian is

$$\mathbf{H} = \nabla^2 f(\boldsymbol{\theta}) \in \mathbb{R}^{W \times W}.$$

Explicit construction of  $\mathbf{H}$  is infeasible for modern neural networks. Instead, second-order methods rely on the *Hessian–vector product*

$$\mathbf{H}\mathbf{u} = \nabla^2 f(\boldsymbol{\theta}) \mathbf{u},$$

for a chosen direction  $\mathbf{u} \in \mathbb{R}^W$ .

A key identity enables efficient computation:

$$\mathbf{H}\mathbf{u} = \nabla_{\boldsymbol{\theta}} \left( \nabla f(\boldsymbol{\theta})^\top \mathbf{u} \right).$$

This identity allows Hessian–vector products to be computed using standard backpropagation without forming  $\mathbf{H}$  explicitly.

### 2.11.1 Worked Example: Scalar Hessian–Vector Product

Consider

$$f(w_1, w_2) = w_1^2 + 3w_1w_2 + 2w_2^2.$$

The gradient is

$$\nabla f = \begin{pmatrix} 2w_1 + 3w_2 \\ 3w_1 + 4w_2 \end{pmatrix},$$

and the Hessian is

$$\nabla^2 f = \begin{pmatrix} 2 & 3 \\ 3 & 4 \end{pmatrix}.$$

For  $\mathbf{u} = (1, -1)^\top$ ,

$$\mathbf{H}\mathbf{u} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}.$$

Using the nabla identity,

$$\nabla f^\top \mathbf{u} = -(w_1 + w_2), \quad \nabla_{\mathbf{w}}(-(w_1 + w_2)) = \begin{pmatrix} -1 \\ -1 \end{pmatrix},$$

confirming the result.

### 2.11.2 Hessian–Vector Products via Backpropagation: Concrete Three-Layer Example

**Forward pass.** The forward pass computes

$$\hat{y} = w_3(w_2(w_1x + b_1) + b_2) + b_3,$$

$$\mathcal{L} = \frac{1}{2}(\hat{y} - y)^2.$$

**First backward pass (gradient computation).** Backpropagation yields the gradient of the loss with respect to all parameters:

$$\nabla f(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial \mathcal{L}}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial b_1} \\ \frac{\partial \mathcal{L}}{\partial w_2} \\ \frac{\partial \mathcal{L}}{\partial b_2} \\ \frac{\partial \mathcal{L}}{\partial w_3} \\ \frac{\partial \mathcal{L}}{\partial b_3} \end{pmatrix} = (\hat{y} - y) \begin{pmatrix} w_3w_2x \\ w_3w_2 \\ w_3a_1 \\ w_3 \\ a_2 \\ 1 \end{pmatrix}.$$

This is the standard gradient computed during ordinary backpropagation.

**Concrete three-layer network usage.** To make the role of the Hessian–vector product explicit in a deep learning setting, consider the three-layer scalar network

$$\begin{aligned} z_1 &= w_1x + b_1, & a_1 &= z_1, \\ z_2 &= w_2a_1 + b_2, & a_2 &= z_2, \\ z_3 &= w_3a_2 + b_3, & \hat{y} &= z_3, & \mathcal{L} &= \frac{1}{2}(\hat{y} - y)^2, \end{aligned}$$

with parameter vector

$$\boldsymbol{\theta} = (w_1, b_1, w_2, b_2, w_3, b_3)^\top.$$

**First backward pass (gradient).** Applying standard backpropagation yields

$$\nabla f(\boldsymbol{\theta}) = (\hat{y} - y) \begin{pmatrix} w_3 w_2 x \\ w_3 w_2 \\ w_3 a_1 \\ w_3 \\ a_2 \\ 1 \end{pmatrix}.$$

This is the gradient used by first-order optimizers such as SGD or Adam.

**Scalar construction for curvature probing.** Given a direction vector

$$\mathbf{u} = (u_{w_1}, u_{b_1}, u_{w_2}, u_{b_2}, u_{w_3}, u_{b_3})^\top,$$

we form the scalar

$$s(\boldsymbol{\theta}) = \nabla f(\boldsymbol{\theta})^\top \mathbf{u} = (\hat{y} - y) (u_{w_1} w_3 w_2 x + u_{b_1} w_3 w_2 + u_{w_2} w_3 a_1 + u_{b_2} w_3 + u_{w_3} a_2 + u_{b_3}).$$

This scalar measures how the gradient of the loss changes when the parameters are perturbed along the direction  $\mathbf{u}$ .

**Second backward pass (Hessian-vector product).** Differentiating  $s(\boldsymbol{\theta})$  with respect to the parameters gives

$$\nabla_{\boldsymbol{\theta}} s(\boldsymbol{\theta}) = \nabla^2 f(\boldsymbol{\theta}) \mathbf{u} = \mathbf{H} \mathbf{u}.$$

Each component of this vector captures the second-order sensitivity of the loss with respect to the corresponding parameter along the chosen direction  $\mathbf{u}$ .

**How this fits into forward and backward propagation.** Operationally, the computation can be summarized as:

- **Forward pass:** compute activations  $(a_1, a_2, \hat{y})$  and the loss  $\mathcal{L}$ .
- **First backward pass:** compute the gradient  $\nabla f(\boldsymbol{\theta})$  using standard backpropagation.
- **Scalar projection:** form  $s(\boldsymbol{\theta}) = \nabla f(\boldsymbol{\theta})^\top \mathbf{u}$ .
- **Second backward pass:** backpropagate through  $s(\boldsymbol{\theta})$  to obtain  $\mathbf{H} \mathbf{u}$ .

**Interpretation in deep learning terms.**

- The first backward pass computes *slopes* of the loss.
- The second backward pass computes how those slopes *change* along a specific direction in parameter space.
- The vector  $\mathbf{u}$  selects which direction of curvature is being examined.
- No Hessian matrix is formed; curvature is accessed implicitly through backpropagation.

**Practical meaning.** In large neural networks, the same mechanism applies when  $\boldsymbol{\theta}$  contains millions of parameters. Hessian-vector products computed in this manner are used inside iterative solvers such as conjugate gradient to obtain curvature-aware update directions, analyze sharpness of minima, and enable Hessian-free optimization—all while preserving the standard forward and backward propagation structure used in deep learning frameworks.

**Role of the second backward pass: clarification.** The purpose of the second backward pass depends on how the Hessian–vector product is used. The same mathematical operation serves two distinct roles in deep learning.

- **Standard first-order training (SGD, Adam).**
  - Training consists of one forward pass and one backward pass per batch.
  - The backward pass computes the gradient  $\nabla f(\boldsymbol{\theta})$ .
  - Parameter updates depend only on first-order information.
  - No Hessian–vector products are required.
  - In this setting, a second backward pass is unnecessary.
- **Second backward pass used for analysis or diagnostics.**
  - The second backward pass computes  $\mathbf{H}\mathbf{u} = \nabla^2 f(\boldsymbol{\theta}) \mathbf{u}$ .
  - The direction  $\mathbf{u}$  is chosen externally (e.g., the gradient direction or a random vector).
  - Quantities such as  $\mathbf{u}^\top \mathbf{H}\mathbf{u}$  are used to assess curvature, sharpness, or stability.
  - Training updates are not modified.
  - The second pass is used only to monitor or analyze convergence behavior.
- **Second backward pass used for optimization.**
  - In second-order or Hessian-free methods, the second backward pass directly influences parameter updates.
  - Optimization algorithms such as conjugate gradient solve linear systems of the form
$$\mathbf{H}\mathbf{p} = -\nabla f(\boldsymbol{\theta})$$
using repeated Hessian–vector products.
  - Each Hessian–vector product requires a second backward pass.
  - The resulting direction  $\mathbf{p}$  determines how the model parameters are updated.
  - In this case, the second backward pass is part of the training algorithm itself.
- **Why the second backward pass is not always used.**
  - Hessian–vector products are computationally more expensive than gradients.
  - Stochastic noise in mini-batch training complicates curvature estimation.
  - First-order methods often provide sufficient performance at lower computational cost.
  - As a result, curvature-based methods are applied selectively.

### Summary.

- One backward pass computes the slope of the loss.
- A second backward pass computes how that slope changes.
- The second backward pass may be used for analysis only or as an integral part of optimization.
- Whether it affects training depends on how the Hessian–vector product is used.

**When Hessian-vector products are used in practice.** The use of Hessian-vector products depends on the choice of optimization strategy.

- **Typical deep learning training.**

- Most deep learning models are trained using first-order optimizers such as SGD, Adam, or RMSProp.
- These methods require only one forward pass and one backward pass per batch.
- Curvature information is not explicitly computed.
- Hessian-vector products are therefore *not* used during standard training.

- **Second-order and curvature-aware optimization.**

- When second-order information is desired, Hessian-vector products become essential.
- Methods such as Hessian-free optimization, Newton-CG, and other curvature-aware algorithms rely on repeated evaluations of  $\mathbf{H}\mathbf{u}$ .
- The second backward pass is used to compute how gradients change along selected directions.
- This additional information is used to construct improved update directions for the model parameters.

- **Practical implication.**

- In most applications, first-order methods are preferred due to their simplicity and computational efficiency.
- Hessian-vector products are used selectively, either for advanced optimization or for analyzing the geometry of the loss landscape.
- When a second-order optimizer is employed, the additional backward pass becomes a necessary and integral part of the training procedure.

**Key takeaway.**

- Standard training does not require Hessian-vector products.
- Hessian-vector products become useful and necessary when second-order optimization methods are employed.

**Role and selection of the direction vector in second-order optimization.** Second-order optimization methods do not operate on the Hessian matrix  $\mathbf{H}$  directly. Instead, they rely on Hessian-vector products of the form  $\mathbf{H}\mathbf{u}$ , which require the specification of a direction vector  $\mathbf{u}$  in parameter space.

- **Why a direction vector is required.**

- The Hessian  $\mathbf{H} \in \mathbb{R}^{W \times W}$  is too large to form or store explicitly in deep learning models.
- Second-order methods therefore interact with the Hessian only through its action on vectors.
- The vector  $\mathbf{u}$  specifies the direction along which curvature information is extracted.

- **Direction vectors are not arbitrary parameters.**

- The vector  $\mathbf{u}$  is not learned and is not part of the model.

- It is generated algorithmically by the optimization method.
- Its purpose is to probe curvature or construct update directions.

- **Typical choices of the direction vector  $\mathbf{u}$ .**

- *Gradient direction:*

$$\mathbf{u} = \nabla f(\boldsymbol{\theta}),$$

used to measure curvature along the descent direction and to assess sharpness or stability.

- *Conjugate gradient search directions:*

$$\mathbf{u} = \mathbf{p}_k,$$

where  $\mathbf{p}_k$  is the current search direction generated by the conjugate gradient algorithm when solving

$$\mathbf{H}\mathbf{p} = -\nabla f(\boldsymbol{\theta}).$$

- *Random probe directions:* used to estimate average curvature or spectral properties of the Hessian.

- **How the direction vector affects optimization.**

- The quality of the second-order update depends on the choice of  $\mathbf{u}$ .
- In Newton–CG and Hessian-free optimization, the direction vectors are chosen adaptively to approximate the Newton step.
- Poorly chosen directions yield poor curvature estimates and ineffective updates.

- **Connection to the second backward pass.**

- Each chosen direction  $\mathbf{u}$  triggers a Hessian–vector product  $\mathbf{H}\mathbf{u}$ .
- Computing  $\mathbf{H}\mathbf{u}$  requires a second backward pass through the scalar  $\nabla f(\boldsymbol{\theta})^\top \mathbf{u}$ .
- Multiple directions may be evaluated during a single optimization step.

### Key takeaway.

- Second-order optimization requires both gradient information and carefully chosen direction vectors.
- Hessian–vector products provide curvature information only along those directions.
- The effectiveness of second-order methods depends critically on how the direction vectors are defined and updated.

**Where Newton–CG is used in practice.** Newton–Conjugate Gradient (Newton–CG) methods appear in deep learning and scientific computing primarily in settings where curvature information is valuable and the optimization problem is sufficiently structured.

- **Hessian-free neural network training.**

- Newton–CG is the core optimization engine in Hessian-free training methods.
- The Newton update direction  $\mathbf{p}$  is obtained by approximately solving

$$\mathbf{H}\mathbf{p} = -\nabla f(\boldsymbol{\theta})$$

using conjugate gradient iterations.

- Each CG iteration requires Hessian–vector products  $\mathbf{H}\mathbf{u}$  rather than the full Hessian.
- **Example applications:**
  - \* training deep autoencoders for dimensionality reduction,
  - \* optimization of recurrent neural networks (RNNs) and LSTMs suffering from vanishing or exploding gradients,
  - \* training very deep feedforward networks where SGD converges slowly due to ill-conditioned curvature.
- **Large-scale scientific and engineering optimization.**
  - Newton–CG is widely used in physics-based optimization problems involving smooth, differentiable objective functions.
  - **Example applications:**
    - \* inverse problems in imaging (e.g., tomography and deconvolution),
    - \* parameter estimation in PDE-constrained optimization,
    - \* aerodynamic shape optimization and structural mechanics,
    - \* fluid dynamics and climate modeling simulations.
  - In these problems, curvature information dramatically reduces the number of required optimization iterations.
- **Classical machine learning models.**
  - Newton–CG is commonly used when the loss function is convex or nearly convex.
  - **Example applications:**
    - \* logistic regression with large feature sets,
    - \* generalized linear models (GLMs),
    - \* kernel-based learning methods with smooth regularization.
  - In these settings, Newton–CG often converges in far fewer iterations than gradient descent.
- **Gauss–Newton and Fisher-based methods in deep learning.**
  - Newton–CG is frequently applied to Gauss–Newton or Fisher information matrix approximations of the Hessian.
  - These approximations are positive semi-definite, making CG numerically stable.
  - **Example applications:**
    - \* training neural networks with squared-error or likelihood-based losses,
    - \* natural gradient methods and trust-region approaches,
    - \* curvature-aware fine-tuning of pretrained deep networks.
- **Why Newton–CG is not the default in deep learning.**
  - Each Newton step requires multiple Hessian–vector products, increasing computational cost.
  - Mini-batch stochasticity introduces noise into curvature estimates.
  - First-order methods (SGD, Adam) often achieve competitive performance with significantly lower overhead.
  - As a result, Newton–CG is used selectively rather than as a default optimizer.

**Summary.**

- Newton–CG is most useful when curvature information significantly improves convergence.
- It is widely used in Hessian-free neural network training and large-scale scientific and engineering optimization.
- In deep learning, it is typically applied with Gauss–Newton or Fisher approximations rather than the exact Hessian.
- Its limited use is due to computational cost and stochastic noise, not lack of effectiveness.

**2.11.3 Role of the Second Backward Pass**

- **First-order training (SGD, Adam):** one forward pass and one backward pass compute gradients.
- **Analysis and diagnostics:** a second backward pass computes curvature quantities such as  $\mathbf{u}^\top \mathbf{H} \mathbf{u}$  without affecting updates.
- **Second-order optimization:** Hessian–vector products directly determine update directions in methods such as Newton–CG and Hessian-free optimization.

The second backward pass is therefore optional and task-dependent.

**2.11.4 Direction Vectors in Second-Order Optimization**

The vector  $\mathbf{u}$  is not a model parameter. It is generated by the optimization algorithm.

Typical choices include:

- the gradient direction  $\mathbf{u} = \nabla f(\boldsymbol{\theta})$ ,
- conjugate gradient search directions,
- random probe directions for curvature estimation.

Each choice extracts curvature information along a specific direction.

**2.11.5 Newton–CG and Hessian-Free Methods**

Newton–CG methods solve

$$\mathbf{H} \mathbf{p} = -\nabla f(\boldsymbol{\theta})$$

using conjugate gradient iterations. Each iteration requires a Hessian–vector product.

These methods appear in:

- Hessian-free neural network training,
- RNN and LSTM optimization,
- inverse problems and PDE-constrained optimization,
- Gauss–Newton and Fisher-based approximations.

They are used selectively due to computational cost and stochastic noise.

Table 2.10: Comparison of PyTorch behavior during training and inference

Aspect	Training mode (autograd enabled)	Inference mode (autograd disabled)
Forward pass behavior	Records the full computational graph of tensor operations	Performs forward computation only
Intermediate activations	Stored in memory for backpropagation	Not stored
Gradient tracking	Tracks tensors with <code>requires_grad=True</code>	No gradient tracking
Operation metadata	Saved to enable chain rule application during backward pass	Not saved
Backward pass	Required to compute gradients of the loss	Not performed
Gradient storage	Gradients accumulated in <code>.grad</code> fields of parameters	No gradient buffers allocated
Higher-order derivatives	Supported when explicitly requested (e.g., Hessian-vector products)	Not available
Memory usage	Higher due to graph and activation storage	Significantly reduced
Runtime performance	Slower due to bookkeeping overhead	Faster and more efficient
Typical use case	Model training and optimization	Evaluation, validation, and deployment
How mode is enabled	Default behavior when gradients are required	Using <code>torch.no_grad()</code> and <code>model.eval()</code>

## 2.12 Training vs. Inference in Automatic Differentiation Frameworks

**Key takeaway.** During training, PyTorch retains intermediate computations and constructs a computational graph to enable automatic differentiation. During inference, this bookkeeping is disabled to reduce memory usage and improve execution speed, while preserving the same forward computation.

**Clarification: meaning of the computational graph.** In this context, the term *graph* refers to PyTorch’s dynamic computational graph constructed during the forward pass when automatic differentiation is enabled. This graph is a directed acyclic graph whose nodes correspond to tensor operations and whose edges represent data dependencies between tensors. Each node stores the local derivative information required to apply the chain rule during backpropagation, along with references to its parent tensors. During the backward pass, PyTorch traverses this graph in reverse to compute gradients of the loss with respect to model parameters.

## Chapter 3

# Jacobians in Backpropagation and Second-Order Analysis

### 3.1 Overview

Backpropagation is often introduced algorithmically as a sequence of local gradient computations propagated backward through a network. At a mathematical level, however, backpropagation is an exact application of the multivariable chain rule and can be understood as repeated multiplication by Jacobian transposes.

This chapter presents a unified Jacobian-based interpretation of backpropagation, explains why explicit Jacobian matrices are rarely constructed in practice, and illustrates the concepts through both symbolic and numerical examples. The connection to second-order methods and Hessian–vector products is also clarified.

### 3.2 Jacobian View of the Chain Rule

Consider a feedforward network expressed as a composition of functions:

$$\mathbf{a}_0 = \mathbf{x}, \quad \mathbf{a}_k = f_k(\mathbf{a}_{k-1}), \quad k = 1, \dots, L,$$

with scalar loss

$$\mathcal{L} = \ell(\mathbf{a}_L).$$

The Jacobian of layer  $k$  with respect to its input is

$$\mathbf{J}_k = \frac{\partial \mathbf{a}_k}{\partial \mathbf{a}_{k-1}}.$$

Applying the multivariable chain rule yields

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_0} = \mathbf{J}_1^\top \mathbf{J}_2^\top \cdots \mathbf{J}_L^\top \frac{\partial \mathcal{L}}{\partial \mathbf{a}_L}.$$

Thus, backpropagation is mathematically equivalent to multiplying the upstream gradient by a sequence of Jacobian transposes.

### 3.3 Why Jacobians Are Not Formed Explicitly

If layer  $k$  maps  $\mathbb{R}^{d_{k-1}} \rightarrow \mathbb{R}^{d_k}$ , then

$$\mathbf{J}_k \in \mathbb{R}^{d_k \times d_{k-1}}.$$

In modern neural networks,  $d_k$  and  $d_{k-1}$  may be millions, making explicit construction of  $\mathbf{J}_k$  infeasible.

Backpropagation requires only the product

$$\mathbf{J}_k^\top \mathbf{v},$$

where  $\mathbf{v}$  is the upstream gradient. Automatic differentiation frameworks compute this product implicitly using local derivatives and the chain rule, avoiding explicit Jacobian storage.

This operation is known as a **vector–Jacobian product (VJP)**.

### 3.4 Three-Layer Network: Matrix Form

Consider a three-layer feedforward network:

$$\begin{aligned} \mathbf{a}_0 &= \mathbf{x}, \\ \mathbf{z}_1 &= \mathbf{W}_1 \mathbf{a}_0 + \mathbf{b}_1, \quad \mathbf{a}_1 = \phi(\mathbf{z}_1), \\ \mathbf{z}_2 &= \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2, \quad \mathbf{a}_2 = \psi(\mathbf{z}_2), \\ z_3 &= \mathbf{W}_3 \mathbf{a}_2 + b_3, \quad \hat{y} = \sigma(z_3), \end{aligned}$$

with scalar loss

$$\mathcal{L} = \frac{1}{2}(\hat{y} - y)^2.$$

The Jacobians of each layer with respect to its input are:

$$\begin{aligned} \mathbf{J}_1 &= \text{diag}(\phi'(\mathbf{z}_1)) \mathbf{W}_1, \\ \mathbf{J}_2 &= \text{diag}(\psi'(\mathbf{z}_2)) \mathbf{W}_2, \\ \mathbf{J}_3 &= \sigma'(z_3) \mathbf{W}_3. \end{aligned}$$

Backpropagation computes:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_2} = (\hat{y} - y) \mathbf{J}_3, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{a}_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}_2} \mathbf{J}_2, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{a}_0} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}_1} \mathbf{J}_1.$$

### 3.5 Numerical Example with Explicit Jacobians (with activation and loss derivatives)

We use a three-layer network with three different nonlinearities: Layer 1 uses ReLU, Layer 2 uses tanh, and the output layer uses a sigmoid. The loss is mean squared error (single sample).

#### 3.5.1 Network definition

Let  $d_0 = d_1 = d_2 = 2$  and define

$$\begin{aligned} \mathbf{a}_0 &= \mathbf{x}, \quad \mathbf{z}_1 = \mathbf{W}_1 \mathbf{a}_0 + \mathbf{b}_1, \quad \mathbf{a}_1 = \text{ReLU}(\mathbf{z}_1), \\ \mathbf{z}_2 &= \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2, \quad \mathbf{a}_2 = \tanh(\mathbf{z}_2), \\ z_3 &= \mathbf{W}_3 \mathbf{a}_2 + b_3, \quad \hat{y} = \sigma(z_3), \quad \mathcal{L} = \frac{1}{2}(\hat{y} - y)^2. \end{aligned}$$

We choose

$$\begin{aligned} \mathbf{x} &= \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad \mathbf{W}_1 = \begin{pmatrix} 1 & 2 \\ -0.5 & 1.5 \end{pmatrix}, \quad \mathbf{b}_1 = \begin{pmatrix} 0.1 \\ -0.2 \end{pmatrix}, \\ \mathbf{W}_2 &= \begin{pmatrix} 0.7 & -1.2 \\ 1.1 & 0.3 \end{pmatrix}, \quad \mathbf{b}_2 = \begin{pmatrix} 0.05 \\ -0.1 \end{pmatrix}, \quad \mathbf{W}_3 = \begin{pmatrix} 1.4 & -0.8 \end{pmatrix}, \quad b_3 = 0.2, \quad y = 1. \end{aligned}$$

### 3.5.2 Activation derivatives and loss derivative

The element-wise derivatives used in the Jacobians are:

$$\frac{d}{dz}\text{ReLU}(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}, \quad \frac{d}{dz}\tanh(z) = 1 - \tanh^2(z), \quad \sigma'(z) = \sigma(z)(1 - \sigma(z)).$$

The loss derivative with respect to the prediction is

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \hat{y} - y.$$

### 3.5.3 Forward pass (numerical)

**Layer 1 (ReLU).**

$$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 = \begin{pmatrix} 1 & 2 \\ -0.5 & 1.5 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} + \begin{pmatrix} 0.1 \\ -0.2 \end{pmatrix} = \begin{pmatrix} 4.1 \\ 0.3 \end{pmatrix}.$$

$$\mathbf{a}_1 = \text{ReLU}(\mathbf{z}_1) = \begin{pmatrix} 4.1 \\ 0.3 \end{pmatrix} \quad (\text{both entries are positive}).$$

**Layer 2 (tanh).**

$$\mathbf{z}_2 = \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2 = \begin{pmatrix} 0.7 & -1.2 \\ 1.1 & 0.3 \end{pmatrix} \begin{pmatrix} 4.1 \\ 0.3 \end{pmatrix} + \begin{pmatrix} 0.05 \\ -0.1 \end{pmatrix} = \begin{pmatrix} 2.56 \\ 4.50 \end{pmatrix}.$$

$$\mathbf{a}_2 = \tanh(\mathbf{z}_2) = \begin{pmatrix} \tanh(2.56) \\ \tanh(4.50) \end{pmatrix} = \begin{pmatrix} 0.988119 \\ 0.999753 \end{pmatrix}.$$

**Output layer (sigmoid).**

$$z_3 = \mathbf{W}_3 \mathbf{a}_2 + b_3 = \begin{pmatrix} 1.4 & -0.8 \end{pmatrix} \begin{pmatrix} 0.988119 \\ 0.999753 \end{pmatrix} + 0.2 = 0.783564.$$

$$\hat{y} = \sigma(z_3) = \frac{1}{1 + e^{-z_3}} = 0.686448, \quad \mathcal{L} = \frac{1}{2}(\hat{y} - y)^2 = \frac{1}{2}(0.686448 - 1)^2.$$

### 3.5.4 Local derivatives (numerical)

**Loss derivative.**

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \hat{y} - y = 0.686448 - 1 = -0.313552.$$

**Sigmoid derivative.**

$$\sigma'(z_3) = \hat{y}(1 - \hat{y}) = 0.686448(1 - 0.686448) = 0.215271.$$

**tanh derivative (element-wise).**

$$\tanh'(\mathbf{z}_2) = 1 - \mathbf{a}_2^{\odot 2} = \begin{pmatrix} 1 - 0.988119^2 \\ 1 - 0.999753^2 \end{pmatrix} = \begin{pmatrix} 0.023620 \\ 0.000494 \end{pmatrix}.$$

**ReLU derivative (element-wise).** Since  $\mathbf{z}_1 = (4.1, 0.3)^\top$  is strictly positive component-wise,

$$\text{ReLU}'(\mathbf{z}_1) = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

### 3.5.5 Explicit Jacobians

**Jacobian of Layer 1:**  $\mathbf{a}_1 = \text{ReLU}(\mathbf{z}_1)$ ,  $\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$ .

$$\mathbf{J}_1 = \frac{\partial \mathbf{a}_1}{\partial \mathbf{x}} = \text{diag}(\text{ReLU}'(\mathbf{z}_1)) \mathbf{W}_1 = \mathbf{I} \mathbf{W}_1 = \mathbf{W}_1.$$

**Jacobian of Layer 2:**  $\mathbf{a}_2 = \tanh(\mathbf{z}_2)$ ,  $\mathbf{z}_2 = \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2$ .

$$\mathbf{J}_2 = \frac{\partial \mathbf{a}_2}{\partial \mathbf{a}_1} = \text{diag}(\tanh'(\mathbf{z}_2)) \mathbf{W}_2 = \begin{pmatrix} 0.023620 & 0 \\ 0 & 0.000494 \end{pmatrix} \begin{pmatrix} 0.7 & -1.2 \\ 1.1 & 0.3 \end{pmatrix} = \begin{pmatrix} 0.016535 & -0.028345 \\ 0.000543 & 0.000148 \end{pmatrix}.$$

**Jacobian of Output:**  $\hat{y} = \sigma(z_3)$ ,  $z_3 = \mathbf{W}_3 \mathbf{a}_2 + b_3$ .

$$\mathbf{J}_3 = \frac{\partial \hat{y}}{\partial \mathbf{a}_2} = \sigma'(z_3) \mathbf{W}_3 = 0.215271 \begin{pmatrix} 1.4 & -0.8 \end{pmatrix} = \begin{pmatrix} 0.301332 & -0.172190 \end{pmatrix}.$$

### 3.5.6 Backpropagation via Jacobian products (numerical)

Backpropagation is a sequence of vector–Jacobian products:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_2} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{a}_2} = (\hat{y} - y) \mathbf{J}_3.$$

Numerically,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_2} = (-0.313552) \begin{pmatrix} 0.301332 & -0.172190 \end{pmatrix} = \begin{pmatrix} -0.094506 & 0.053988 \end{pmatrix}.$$

Next,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}_2} \mathbf{J}_2,$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_1} = \begin{pmatrix} -0.094506 & 0.053988 \end{pmatrix} \begin{pmatrix} 0.016535 & -0.028345 \\ 0.000543 & 0.000148 \end{pmatrix} = \begin{pmatrix} -0.001435 & 0.002267 \end{pmatrix}.$$

Finally,

$$\nabla_{\mathbf{x}} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}_1} \mathbf{J}_1 = \frac{\partial \mathcal{L}}{\partial \mathbf{a}_1} \mathbf{W}_1,$$

$$\nabla_{\mathbf{x}} \mathcal{L} = \begin{pmatrix} -0.001435 & 0.002267 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ -0.5 & 1.5 \end{pmatrix} = \begin{pmatrix} -0.002876 & 0.000531 \end{pmatrix}.$$

Equivalently, written as a column vector:

$$\nabla_{\mathbf{x}} \mathcal{L} = \begin{pmatrix} -0.002876 \\ 0.000531 \end{pmatrix}.$$

**Interpretation.**

- The loss derivative  $\partial\mathcal{L}/\partial\hat{y}$  produces the initial upstream signal.
- The output Jacobian  $\mathbf{J}_3$  injects the sigmoid derivative and the linear weights  $\mathbf{W}_3$ .
- The hidden Jacobian  $\mathbf{J}_2$  injects the tanh derivative via a diagonal gate.
- The input Jacobian  $\mathbf{J}_1$  injects the ReLU derivative via another diagonal gate.
- At no point is any full Jacobian of the entire network constructed; only local Jacobians are used in products.

**Note.** Minor differences at the last digit can occur depending on rounding of  $\tanh(\cdot)$  and  $\sigma(\cdot)$ .

### 3.6 Connection to Second-Order Derivatives

The Jacobian framework naturally extends to second-order analysis. Given the gradient  $\nabla f(\boldsymbol{\theta})$ , the Hessian–vector product is defined as

$$\mathbf{H}\mathbf{u} = \nabla_{\boldsymbol{\theta}}(\nabla f(\boldsymbol{\theta})^\top \mathbf{u}).$$

This computation requires:

- one forward pass,
- one backward pass to compute  $\nabla f(\boldsymbol{\theta})$ ,
- formation of the scalar  $\nabla f(\boldsymbol{\theta})^\top \mathbf{u}$ ,
- a second backward pass.

No Hessian matrix is ever formed.

### 3.7 Summary

- Backpropagation is mathematically equivalent to multiplying by Jacobian transposes.
- Automatic differentiation computes vector–Jacobian products implicitly.
- Explicit Jacobian matrices are avoided due to prohibitive size.
- The same machinery enables efficient computation of Hessian–vector products.
- First-order and second-order optimization methods differ only in how many backward passes they require.



## Chapter 4

# Vision Transformer (ViT)

### 4.1 Overview

The Vision Transformer (ViT) is a neural architecture that applies the Transformer model (originally developed for natural language processing) to image understanding tasks. Instead of relying on convolutional operations, ViT represents an image as a sequence of embedded patches and processes this sequence using stacked self-attention layers.

Unlike convolutional neural networks (which impose locality and translation equivariance through kernel design), ViT relies on global self-attention to model long-range spatial relationships. This design choice introduces different scaling behaviors in memory usage and computational cost, which must be carefully considered when ViT is used as a backbone in detection or segmentation frameworks such as Mask R-CNN.

### 4.2 Patch Embedding and Tokenization

Given an input image of size  $H \times W \times C$  (height  $\times$  width  $\times$  channels), ViT first partitions the image into non-overlapping square patches of size  $P \times P$ . Each patch is flattened and projected into a fixed-dimensional embedding space.

The number of patch tokens is given by:

$$N_{\text{patch}} = \frac{H}{P} \cdot \frac{W}{P} \quad (4.1)$$

A special learnable classification token ([CLS]) is prepended to the patch sequence, yielding the total number of tokens:

$$N = N_{\text{patch}} + 1 \quad (4.2)$$

Each token is represented as a vector in  $\mathbb{R}^D$ , where  $D$  is the embedding dimension.

### 4.3 Transformer Depth and Attention Modules

A Vision Transformer consists of a stack of identical Transformer layers. Each layer contains exactly one multi-head self-attention (MHSA) module followed by a feed-forward multilayer perceptron (MLP).

- The number of **attention modules** is equal to the Transformer depth  $L$ .
- Each attention module processes the full token sequence.
- Depth is a user-defined architectural hyperparameter and is independent of token count, embedding dimension, and number of attention heads.

Thus:

$$\text{Number of attention modules} = L \quad (4.3)$$

## 4.4 Multi-Head Self-Attention

Within a single attention module, multi-head self-attention projects the token embeddings into queries, keys, and values, which are then split across multiple attention heads.

Let:

- $D$  be the embedding dimension
- $h$  be the number of attention heads

The per-head dimension is:

$$d_h = \frac{D}{h}, \quad \text{with } D \bmod h = 0 \quad (4.4)$$

For each head, attention is computed as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{\sqrt{d_h}} \right) V \quad (4.5)$$

The attention matrix for each head has size:

$$N \times N \quad (4.6)$$

This quadratic dependence on the number of tokens is the primary source of memory and computational cost in ViT.

## 4.5 Computational Scaling

The dominant computational and memory terms in ViT arise from self-attention and MLP layers.

### 4.5.1 Attention Cost

Self-attention scales as:

$$\mathcal{O}(N^2 \cdot D) \quad (4.7)$$

This quadratic scaling with respect to the number of tokens makes patch size and input resolution critical design parameters.

### 4.5.2 MLP Cost

Each Transformer layer contains an MLP with hidden dimension:

$$D_{\text{MLP}} = r \cdot D \quad (4.8)$$

where  $r$  is the MLP expansion ratio, typically set to 4.

The MLP cost scales as:

$$\mathcal{O}(N \cdot D \cdot D_{\text{MLP}}) \quad (4.9)$$

## 4.6 Parameter Impact Summary

Table 4.1 summarizes the role and impact of key Vision Transformer parameters.

Parameter	Controls	Primary Impact
Patch size $P$	Number of tokens $N$	Quadratic attention memory and compute
Embedding dim $D$	Token feature width	Linear memory, quadratic MLP cost
Depth $L$	Number of attention modules	Linear scaling in memory and compute
Heads $h$	Attention parallelism	Minor overhead, $D \bmod h = 0$
MLP ratio $r$	Feed-forward width	Significant compute increase
Image resolution $H \times W$	Token count	Quadratic attention scaling

Table 4.1: Impact of Vision Transformer architectural parameters.

## 4.7 Numerical Example

Consider an input image of size  $224 \times 224$  with patch size  $P = 16$ .

$$N_{\text{patch}} = (224/16)^2 = 14^2 = 196 \quad (4.10)$$

$$N = 196 + 1 = 197 \quad (4.11)$$

For an embedding dimension  $D = 192$  and  $h = 3$  heads:

$$d_h = 192/3 = 64 \quad (4.12)$$

The attention matrix size per layer has size:

$$h \cdot N^2 = 3 \times 197^2 = 116,427 \quad (4.13)$$

If the patch size is reduced to  $P = 8$ , the number of tokens increases to  $N = 785$ , and the attention matrix size increases to approximately 1.85 million entries per layer, illustrating the quadratic scaling effect.

## 4.8 Implications for Detection and Segmentation

When used as a backbone in detection or segmentation frameworks, ViT must often process images at higher resolutions than standard classification benchmarks. As a result, token count and attention memory become the dominant constraints.

Practical deployments typically rely on:

- Larger patch sizes
- Reduced depth
- Windowed or local attention mechanisms
- Hierarchical token merging

These modifications preserve the representational benefits of self-attention while keeping memory usage within feasible limits.

## 4.9 Summary

Vision Transformers replace convolutional inductive biases with global self-attention over patch tokens. The number of attention modules is determined solely by the Transformer depth, while computational cost is dominated by the quadratic dependence on the number of tokens.

Parameter	Controls	Compute / Memory Scaling	Effect on Representati
MLP ratio $r$	Hidden width of MLP ( $rD$ )	$\mathcal{O}(N \cdot rD^2)$ compute, $\mathcal{O}(N \cdot rD)$ activations	Channel-wise capacity (no

Table 4.2: Impact of the MLP ratio in Vision Transformers.

Understanding these scaling relationships is essential when integrating ViT into large-scale vision systems such as Mask R-CNN.

The MLP ratio determines the width of the channel-wise feed-forward network applied independently to each token, increasing parameter count, compute, and activation memory linearly while leaving token-token interactions unchanged, which are governed exclusively by self-attention. In contrast to self-attention, which mixes information across tokens, the MLP operates only along the feature dimension and therefore controls expressive capacity rather than spatial interaction.

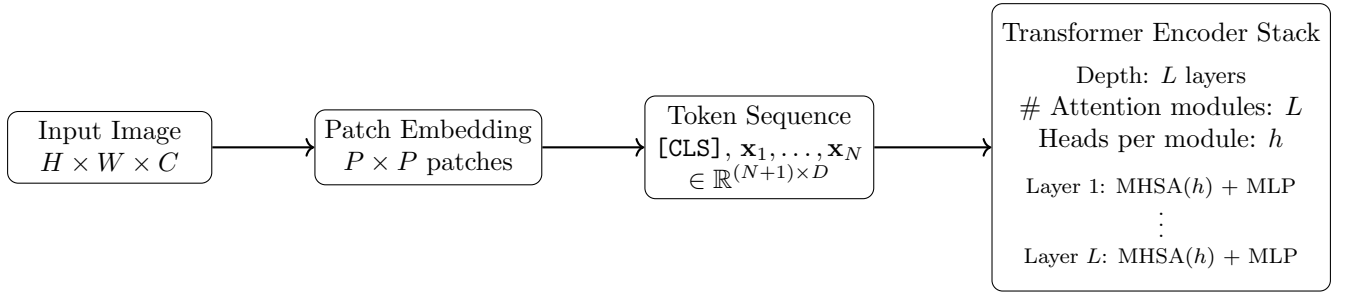


Figure 4.1: High-level architecture of the Vision Transformer (ViT). An input image is split into non-overlapping patches, embedded into a token sequence, and processed by a Transformer encoder stack of depth  $L$ . Each encoder layer contains exactly one multi-head self-attention (MHSA) module (so the number of attention modules is  $L$ ), and each MHSA uses  $h$  attention heads.

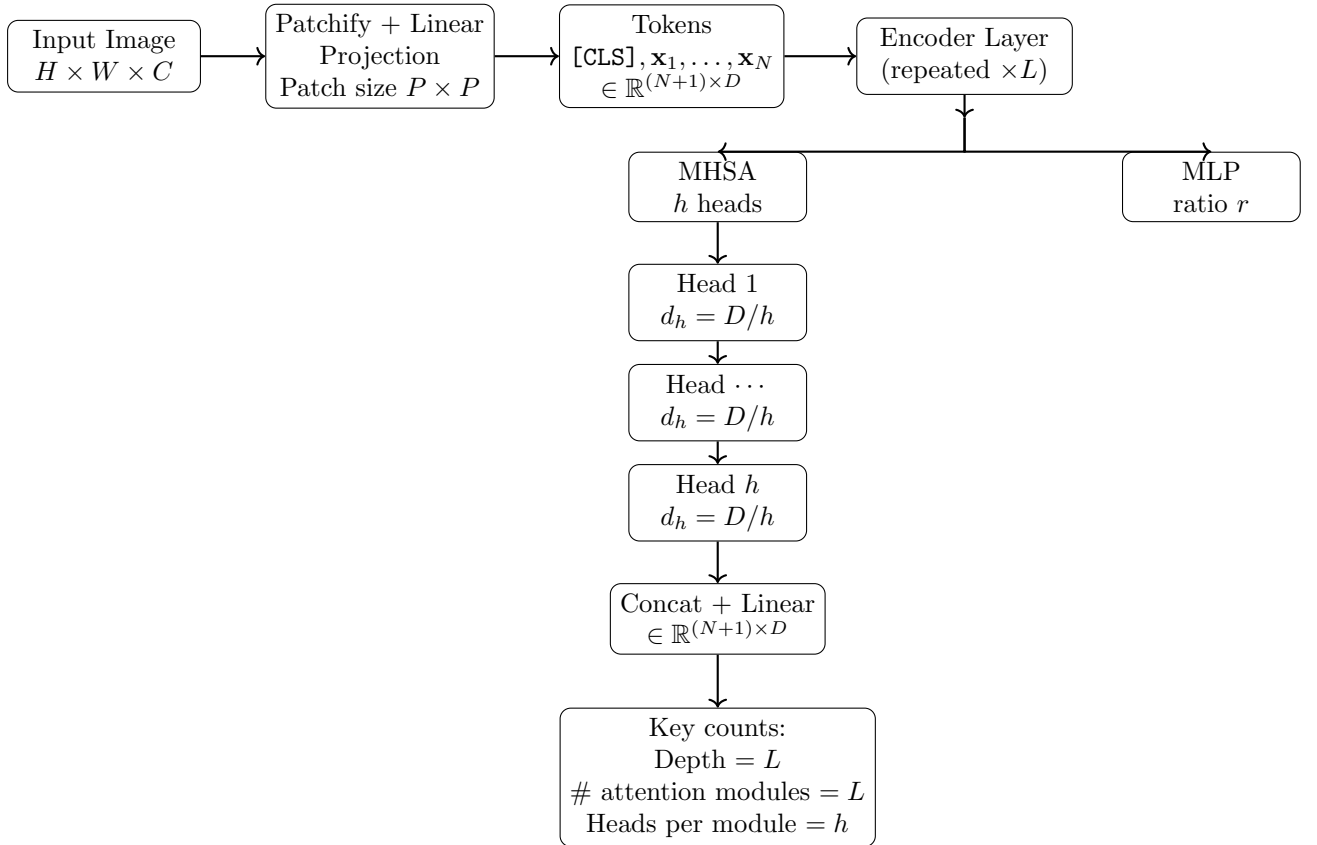


Figure 4.2: Hydra/tree-style view of ViT without overflow. Tokens pass through an encoder layer repeated  $\times L$ . Each layer contains exactly one MHSA (so the number of attention modules is  $L$ ) with  $h$  heads and one MLP block. The MHSA splits into  $h$  heads and merges back via concatenation to dimension  $D$ .

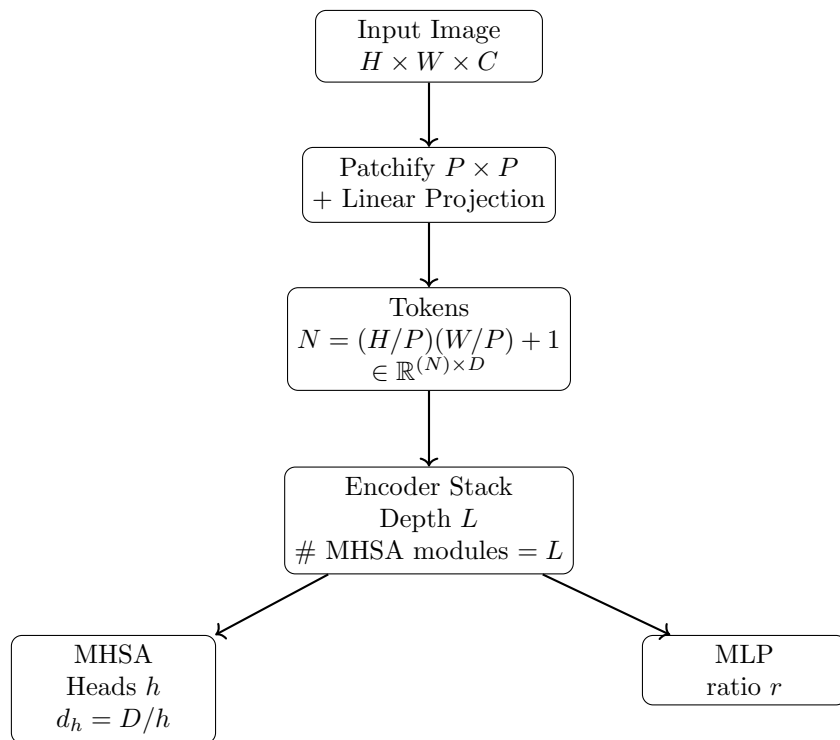


Figure 4.3: Tree-style view of ViT. The encoder stack has depth  $L$ ; each layer contains one MHA (therefore total attention modules =  $L$ ) with  $h$  heads and an MLP block.

# Index

- Activation function, [9](#), [10](#), [15](#)
  - choice, [10](#)
  - derivative, [28](#)
  - ReLU, [11](#)
  - sigmoid, [11](#)
  - tanh, [11](#)
- Activation Functions, [1](#)
  - GELU, [2](#)
  - Leaky ReLU, [2](#)
  - ReLU, [2](#)
  - Sigmoid, [1](#)
  - Softmax, [2](#)
  - Tanh, [1](#)
- Adagrad, [10](#)
- Adam, [10](#)
- Architecture parameters, [34](#)
- Attention
  - cost, [34](#)
- Attention head, [33–35](#)
  - count, [34](#), [36–38](#)
- Attention matrix, [34](#)
  - size, [35](#)
- Attention module, [33](#), [34](#)
  - count, [33](#), [35–38](#)
- Autograd, [26](#)
- Automatic differentiation, [9](#), [14](#), [26](#), [27](#)
- Backbone network, [33](#), [35](#)
- Backpropagation, [1](#), [7](#), [9](#)
  - definition, [9](#)
  - gradient accumulation, [17](#)
  - Jacobian view, [27](#)
  - matrix form, [15](#)
  - single neuron, [11](#)
  - steps, [10](#)
  - summary, [10](#), [31](#)
  - universality, [10](#)
- Backward pass, [9](#)
- Bias, [9](#), [11](#)
- Chain rule, [9–11](#), [27](#)
  - accumulation, [17](#)
- Channels, [33](#)
- [CLS] token, [33](#)
- Computation graph, [11](#), [14](#), [26](#)
  - definition, [26](#)
- Computational complexity, [27](#), [33–36](#)
- Compute, [34](#), [35](#)
- Computer vision, [33](#)
  - systems, [36](#)
  - Transformers, [33](#)
- Convolution, [33](#)
- Convolutional neural network (CNN), [11](#), [33](#)
- Deep neural network (DNN), [11](#)
- Delta (backprop), [15](#)
- Deployment, [35](#)
- Depth, [10](#)
- Depth  $L$ , [33](#), [35](#)
- Derivative
  - activation function, [15](#)
  - ReLU, [12](#)
  - sigmoid, [12](#)
  - tanh, [12](#)
- Direction vector, [24](#)
- Divisibility constraint, [35](#)
- Embedding, [33](#), [34](#)
- Embedding dimension, [33](#), [35](#)
  - per-head, [34](#)
- Embedding dimension  $D$ , [33](#), [34](#)
- Example
  - attention matrix, [35](#)
  - numerical values, [12](#)
  - single neuron, [11](#)
  - three-layer backpropagation, [11](#)
  - token count, [35](#)
- Feed-forward network, [33](#), [36](#)
  - cost, [34](#)
- Feedforward neural network, [9](#)
- Flattening, [33](#)
- Forward Pass, [1](#), [7](#)
- Forward pass, [9](#)
  - data flow, [9](#)
  - definition, [9](#)
  - matrix form, [15](#)
  - memory, [15](#), [16](#)

- single neuron, 11
- summary, 10
- Gradient, 9
  - biases, 16
  - weights, 16
- Gradient descent, 10
- Hadamard product, 15
- Head dimension, 34
- Heads  $h$ , 34, 35
- Hessian-free optimization, 18, 24
- Hessian-vector product, 17, 18, 31
  - example, 18
- Hierarchical transformer, 35
- Hyperparameter, 33, 34
  - typical values, 34
- Image, 33
- Image classification, 35
- Image resolution, 33–35
- Image segmentation, 33, 35
- Inductive bias, 35
  - locality, 33
  - translation equivariance, 33
- Inference mode, 26
- Jacobian matrix, 27
  - definition, 27
  - numerical example, 28
- Jacobian–vector product, 27
- Kernel, 33
- Key (K), 34
- Layer
  - hidden, 9
  - input, 9
  - output, 9
- Linear projection, 33
- Long-range dependency, 33
- Loss function, 9
  - choice, 10
  - derivative, 28
- Loss Functions, 1, 2
  - Cross Entropy, 3
  - cross-entropy, 9, 10
  - Dice Loss, 3
  - Focal Loss, 3
  - Huber Loss, 3
  - IoU Loss, 3
  - KL Divergence, 3
  - MAE, 2
- MSE, 2, 9, 10
- Mask R-CNN, 33, 36
- Matrix calculus, 15
- Matrix–vector product, 15
- Memory complexity, 27
- Memory usage, 26, 33–36
  - attention, 35
- MLP, 33, 34
  - cost, 34, 35
  - width, 35
- MLP ratio, 36
- MLP ratio  $r$ , 34, 35
- Multi-head self-attention (MHSA), 33, 34, 36, 37
- Multilayer perceptron (MLP), 9, 28
- Nabla operator, 17
- Natural language processing (NLP), 33
- Network architecture
  - layers, 9
- Neural network
  - basic, 9
  - training, 9
- Newton–CG, 18, 24
- Notation
  - dimensions, 15
- Object detection, 33, 35
- Optimizer, 10
  - choice, 10
- Optimizers, 1, 4
  - Adagrad, 5
  - Adam, 5
  - AdamW, 5
  - LAMB, 5
  - Momentum, 5
  - RMSProp, 5
  - SGD, 4
- Outer product, 15, 16
- Parallelism, 35
- Patch, 33
  - non-overlapping, 33
- Patch embedding, 33, 36
- Patch size, 33–35
- Product types, 16
- PyTorch, 26
- Quadratic complexity, 34, 35
- Query (Q), 34
- Recurrent neural network (RNN), 11

- ReLU, 9, 10
  - derivative, 15
- RMSProp, 10
- Saved tensors, 15
- Scaling laws, 33–36
- Second backward pass, 18, 24
- Second-order derivatives, 17
- Second-order optimization, 27
- Self-attention, 33–35
  - cost, 34
  - global, 33, 35
  - local, 35
  - quadratic scaling, 34, 35
- Sequence
  - prepend token, 33
- Sequence model, 33
- Sigmoid, 9, 10
  - derivative, 15
- Stacked layers, 33
- Stochastic gradient descent (SGD), 10
- tanh, 9, 10
  - derivative, 15
- Target, 9
- Token, 33, 36
  - classification token, 33
  - count, 33–35
  - embedding, 35
  - merging, 35
  - patch token, 33, 35
  - sequence, 33
- Tokenization, 33
- Training mode, 26
- Transformer, 33
  - depth, 33, 35–38
  - encoder, 36
  - layer, 33, 34
  - overview, 33
- Value (V), 34
- Vector–Jacobian product, 27
- Vision Transformer (ViT), 33, 35
  - architecture, 36–38
  - parameters, 34
  - summary, 35
- ViT, *see* Vision Transformer (ViT)
- Weight and Bias Updates, 1, 7
- Weighted sum, 9
- Weights, 11
- Width, 10
- Windowed attention, 35