

# **Modern Deep Learning Foundations**

Course Notes

Dr. Barak Or

August 2025

Artificial**Gate**

# Contents

<b>Introduction</b>	<b>4</b>
<b>1 Machine Learning vs. Deep Learning</b>	<b>6</b>
<b>2 What Is a Neural Network</b>	<b>11</b>
<b>3 Loss Function, Backpropagation, Optimization</b>	<b>17</b>
<b>4 How Does Training Actually Work?</b>	<b>24</b>
<b>5 Performance Evaluation Metrics</b>	<b>31</b>
<b>6 Overfitting and Regularization</b>	<b>39</b>
<b>7 Why Do We Need Convolution?</b>	<b>45</b>
<b>8 How Does a CNN Work?</b>	<b>52</b>
<b>9 Sequences and Time: RNN, GRU, and LSTM</b>	<b>58</b>
<b>10 Autoencoders for Dimensionality Reduction</b>	<b>64</b>
<b>11 Self-Attention and the Transformer Principle</b>	<b>70</b>
<b>12 Normalization and Initialization</b>	<b>75</b>

<b>13 Data Augmentation</b>	<b>80</b>
<b>14 Advanced Optimization</b>	<b>85</b>
<b>15 Explainability: Understanding Model Decisions</b>	<b>91</b>
<b>16 TensorFlow vs. PyTorch</b>	<b>98</b>
<b>17 Working Effectively with Google Colab</b>	<b>104</b>
<b>18 Mixed Precision Training</b>	<b>111</b>
<b>19 Transfer Learning and Fine-Tuning</b>	<b>116</b>
<b>20 Saving, Loading, and Versioning Models</b>	<b>122</b>
<b>21 Basic Industrial Deployment</b>	<b>128</b>
<b>22 Advancing into Specialized Domains</b>	<b>134</b>
<b>23 Roadmap for the Industrial DL Engineer</b>	<b>139</b>

# Introduction

## Welcome.

This book is the result of years of hands-on work in deep learning - training engineers, advising companies, and deploying models that operate at scale. It's based on the foundational deep learning course I teach through ArtificialGate, and it's been carefully crafted to take you from zero to deployment.

My name is Dr. Barak Or. I hold a PhD in machine learning for navigation systems. Over the past decade, I've founded startups and integrated deep learning into real-world systems. More information on my personal website: [www.barakor.com](http://www.barakor.com).

This book is not a research survey or a theoretical treatment. It's a foundation for engineers. It is designed to give you everything you need to:

- Understand how deep learning works from first principles
- Write, train, and evaluate models in PyTorch
- Save and deploy models in real production-like environments
- Choose the right tools, formats, and practices for modern AI work

You'll find full examples, end-to-end pipelines, deployment walkthroughs, and engineering checklists - all written in clean, accessible language, without cutting corners.

Whether you're an engineer new to AI, or a data scientist moving into deployment, this book is designed to help you build confidently.

**Who is this book for?** It assumes you know Python and have some background in math - but not necessarily in machine learning. It's designed for those who want to move beyond notebooks and into systems - building tools that serve users and last beyond the experiment.

**How to use it:** This book serves as the official course notes for the interactive program on ArtificialGate.ai. Alongside the text, the platform provides a complete learning environment: Google Colab notebooks with exercises and solutions, video-based lessons, and guided walkthroughs. The course is available in many languages and is designed to support learners worldwide across domains and backgrounds.

The lessons here are designed to be read sequentially, but each also stands on its own. You can run the code, follow along in Colab, or simply absorb the concepts and return to the code when you're ready. Each chapter builds intuition, introduces a key tool or practice, and brings you one step closer to deploying deep learning systems with confidence.

# Lesson 1

## Machine Learning vs. Deep Learning

### Introduction

Before we understand what a neural network looks like, or what gradient descent means, it is important to pause and build a foundation: what exactly is machine learning? What distinguishes it from deep learning? And why is this distinction critical for engineers, researchers, and developers?

### Machine Learning: Basic Principles

Rather than writing explicit rules such as: “If an email says you won a million dollars - mark it as spam,” we provide the model with thousands of emails labeled as either legitimate or spam, and the model learns on its own what separates the two categories.

Every such model undergoes a process known as training. In supervised learning, we provide pairs of inputs and labels - for example, a picture of a cat labeled as “1” and a picture of a dog labeled as “0.” The model makes a prediction based on the input alone, compares the prediction to the actual label, calculates the error, and updates itself so that the error decreases. This happens thousands of times - until

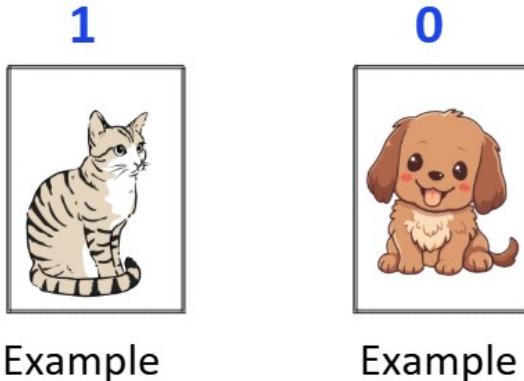


Figure 1.1: Images of a cat and a dog labeled “1” and “0” respectively.

the model learns the relationship between the input (the image) and the desired output. In other words, the goal of training is to reduce the discrepancy between the true label and the model’s prediction.

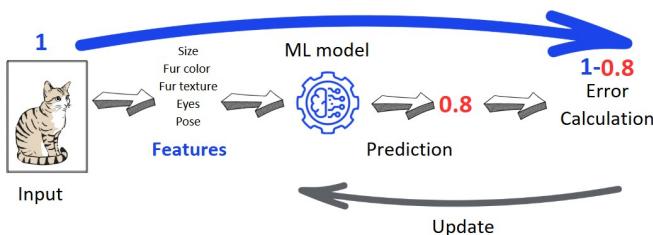


Figure 1.2: Illustration of the training process: prediction, error calculation, and update loop.

During training, the model learns to associate features - such as height, color, or texture - with the correct outcome. In classical machine learning, we must define these features ourselves: decide what to feed into the model, in what format, and how to scale or balance the values.

## The Role of Deep Learning

This is exactly where deep learning comes in. Deep learning is not a single technique, but a family of methods based on neural networks. Here, the model receives input, passes it through many layers of transformation, and builds internal representations - without requiring manual feature design. If the input is an image, for example, the model receives the full set of pixels rather than pre-engineered features.

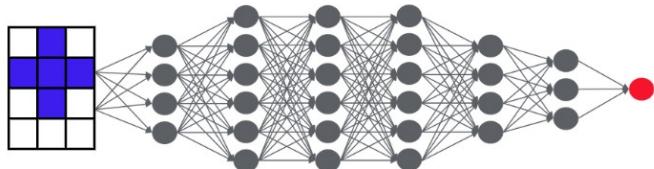


Figure 1.3: A deep neural network that takes raw image pixels as input.

## Case Reflection: Feature Discovery in Practice

In one applied project, the task was to classify product images based on perceived attributes such as luxury, brightness, and cleanliness. The initial approach, using conventional machine learning, relied on structured metadata features-price, category, and user ratings. While this produced reasonable results, performance was limited by the expressiveness of manually engineered features.

Transitioning to a deep learning approach, we trained models directly on the raw image data without predefined features. The model was able to learn internal representations that captured subtle visual patterns that could not be easily defined manually. This underscored a key advantage of deep learning: the ability to discover and extract

relevant features automatically, without explicit human specification.

## Hierarchical Representation Learning

Deep learning is a process of hierarchical representation learning: layer by layer, the information is processed and represented at increasing levels of abstraction - down to the level of individual pixels in the image.

This is what makes deep learning so powerful, especially with unstructured data. Since 2012, beginning with the breakthrough of AlexNet, deep learning has transformed fields such as computer vision, speech recognition, translation, sentiment analysis, and content generation.

## When Is Deep Learning the Right Choice?

Deep learning is the right approach when dealing with unstructured data, when there is a large amount of data available, when we want the system to learn representations independently, or when we aim to uncover complex patterns that are difficult to define manually.

## Challenges and Tradeoffs

Despite its significant capabilities, deep learning presents several challenges that must be considered in practical deployment. Training deep networks requires substantial computational resources and often involves long optimization cycles. Moreover, these models tend to exhibit low interpretability: it is often non-trivial to explain why a specific prediction was made, as the internal representations are distributed across many layers and parameters.

This characteristic has led to the common characterization of deep

networks as “black boxes”-systems in which the decision process is difficult to trace. Nevertheless, when trained and validated appropriately, these architectures have demonstrated the ability to extract structure from high-dimensional data and uncover patterns that are not accessible through traditional methods.

## Summary

Deep learning is not simply an advanced version of machine learning. It represents a shift in modeling philosophy. Rather than manually designing features, we build systems that learn to represent data through layered abstraction.

The key differences can be summarized as follows: classical machine learning relies on manual feature design, flat architectures, and structured data. Deep learning performs automatic feature learning, uses deep network structures, and is particularly effective for unstructured data.

In the upcoming lessons, we will explore how to construct such models - starting from the basics of a single artificial neuron and progressing to complete deep networks.

## Lesson 2

### What Is a Neural Network

#### Introduction

In the previous lesson, we explored the difference between machine learning and deep learning. We saw that deep learning allows the model to learn internal representations by itself, without requiring manual feature engineering. In this lesson, we examine the internal structure and fundamental mathematical formulation of artificial neural networks.

#### The Artificial Neuron

At the core of every neural network lies the artificial neuron, also referred to as a perceptron. A neuron receives a vector of input features  $\mathbf{x} \in \mathbb{R}^n$ , multiplies each input by a corresponding weight, sums the results, adds a bias term  $b \in \mathbb{R}$ , and passes the result through an activation function  $\varphi$  to produce the output  $y$ .

Formally, the neuron computes:

$$y = \varphi(\mathbf{w} \cdot \mathbf{x} + b) \quad (2.1)$$

Where:

- $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$  is the input vector.
- $\mathbf{w} = [w_1, w_2, \dots, w_n]^\top$  is the weight vector.
- $\mathbf{w} \cdot \mathbf{x}$  denotes the dot product  $\sum_{i=1}^n w_i x_i$ .
- $b$  is a scalar bias term.
- $\varphi(\cdot)$  is a non-linear activation function applied to the scalar pre-activation value.
- $y$  is the scalar output of the neuron.

While the neuron is a simple function, connecting many such neurons into layers enables the modeling of highly complex functions.

## Activation Functions

The activation function  $\varphi$  introduces non-linearity into the model, which is essential for learning complex and non-linear patterns. If a linear function is used (e.g.,  $\varphi(x) = 3x$ ), the network behaves as a linear transformation regardless of its depth. Non-linear activation functions overcome this limitation.

Two commonly used activation functions are:

- **ReLU (Rectified Linear Unit):**

$$\varphi(x) = \max(0, x) \tag{2.2}$$

ReLU outputs zero for negative inputs and passes positive inputs unchanged. It is simple, computationally efficient, and widely used in modern deep networks.

- **Sigmoid:**

$$\varphi(x) = \frac{1}{1 + e^{-x}} \tag{2.3}$$

Sigmoid squashes inputs into the range  $(0, 1)$ . It was traditionally used in neural networks, especially for binary classification. However, its use has declined due to the *vanishing gradient* problem, where gradients diminish exponentially during backpropagation, making it difficult to update parameters in deep networks.

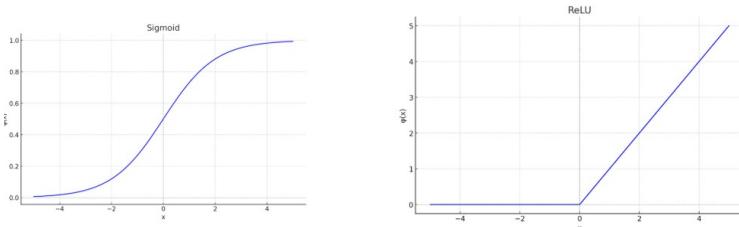


Figure 2.1: Activation functions: Sigmoid (left) and ReLU (right).

As a practical tip: if the training loss is not decreasing, one common reason is the improper use of activation functions-such as using Sigmoid in hidden layers instead of ReLU.

## Layered Network Structures

A neural network is constructed by stacking multiple layers of neurons. Each layer takes the output of the previous layer as its input and produces a new set of activations.

Formally, for layer  $\ell$ , the output is computed as:

$$\mathbf{y}^{(\ell)} = \varphi(\mathbf{W}^{(\ell)}\mathbf{y}^{(\ell-1)} + \mathbf{b}^{(\ell)}) \quad (2.4)$$

Where:

- $\mathbf{y}^{(\ell-1)}$  is the input vector to layer  $\ell$  (i.e., output of layer  $\ell - 1$ ),
- $\mathbf{W}^{(\ell)}$  is the weight matrix for layer  $\ell$ ,

- $\mathbf{b}^{(\ell)}$  is the bias vector for layer  $\ell$ ,
- $\varphi$  is applied element-wise to the result.

Network architecture includes:

- **Input Layer:** Get the raw input data (e.g., image pixels or numerical features).
- **Hidden Layers:** Perform non-linear transformations and learn internal representations.
- **Output Layer:** Produces the final prediction, typically with a task-specific activation function.

Each connection in the network has a trainable weight, and each neuron has a trainable bias. During the training process, these parameters are optimized to minimize a loss function.

## Multilayer Perceptron (MLP)

A Multilayer Perceptron is a class of feedforward neural network where each layer is fully connected to the next. This means that every neuron in a given layer receives input from every neuron in the previous layer.

An MLP is especially effective for problems that are not linearly separable. A classic example is the XOR function, defined as:

- Inputs:  $x_1, x_2 \in \{0, 1\}$
- Output:  $y = 1$  if  $x_1 \neq x_2$ , otherwise  $y = 0$

This function cannot be separated by a linear hyperplane.

To model XOR, we can define an MLP with:

- Input layer of 2 neurons (for  $x_1$  and  $x_2$ ),

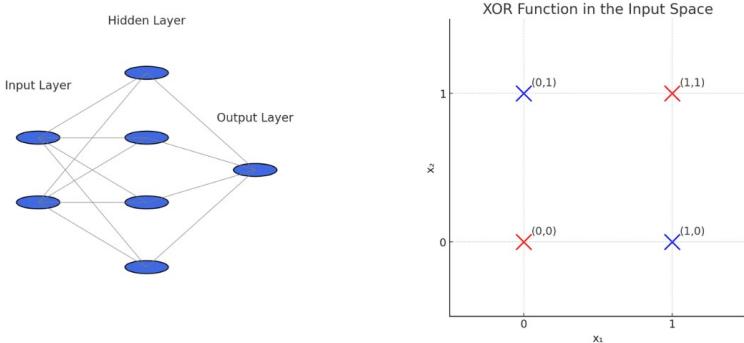


Figure 2.2: XOR classification: simplified neural network (left) and XOR function plot (right).

- Hidden layer with 4 neurons using ReLU,
- Output layer with 1 neuron using Sigmoid.

Without non-linear activations, stacking layers of linear transformations results in a single linear function, which cannot represent XOR. Therefore, non-linearity is essential for modeling complex decision boundaries.

### Universal Approximation in Theory

There is a foundational result in neural network theory known as the *Universal Approximation Theorem*. It states that even a simple feedforward neural network with a single hidden layer can approximate any continuous function on compact input domains, as long as it uses a non-linear activation function (such as ReLU or Sigmoid) and has sufficient neurons. This means that MLPs are powerful enough, in theory, to model almost anything. But in practice, networks with only one large layer are hard to train, inefficient, and prone to overfitting. That's why we use **deep** networks:

stacking layers allows reducing parameter count and generalizing better. The theorem tells us what's possible - deep learning tells us how to get there.

## Summary

A neural network is a parametric function composed of layers of interconnected neurons. Each neuron computes a weighted combination of inputs, adds a bias, and passes the result through a non-linear activation function. The architecture consists of input, hidden, and output layers.

The trainable parameters-weights and biases-are learned through gradient-based optimization. By introducing non-linearity, networks can approximate a wide range of functions, including those that are not linearly separable. Even basic networks like MLPs can solve functions such as XOR, which linear models cannot.

## Lesson 3

# Loss Function, Backpropagation, Optimization

## Introduction

In the previous lesson, we discussed how information flows forward through the layers of a neural network to produce a prediction. However, a network does not begin with knowledge - it must learn from data. This lesson explains how that learning occurs through three essential components: the loss function, the backpropagation algorithm, and the gradient descent optimization method.

## The Loss Function

A neural network learns by minimizing a scalar value known as the **loss**. The loss function quantifies how far the model's prediction is from the actual target label. Its purpose is to measure how "wrong" the prediction is.

Let  $y$  represent the true label (e.g.,  $y \in \{0, 1\}$  in binary classification), and let  $\hat{y}$  represent the predicted probability output by the model. A perfect prediction yields a loss of 0, while incorrect predictions result in a higher loss (see Figure 3.1).

A commonly used loss function for binary classification is the **cross-**

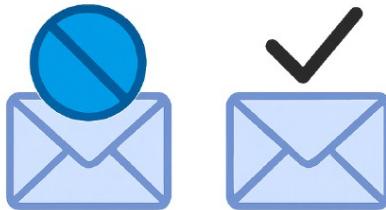


Figure 3.1: Example: one email is spam, the other is not. The network's task is to predict this correctly.

**entropy loss**, defined as:

$$L = -[y \cdot \log(\hat{y}) + (1 - y) \cdot \log(1 - \hat{y})] \quad (3.1)$$

Where:

- $y$  is the true label (0 or 1)
- $\hat{y}$  is the predicted probability
- $L$  is the computed scalar loss

If the prediction is accurate and  $\hat{y}$  is close to  $y$ , the loss is low. If the prediction is far from the true label, the loss is high. The learning objective of the network is to find weights and biases that minimize the average loss across all training examples.

## Training Flow

The complete training process—from input to prediction, loss calculation, and weight update—is summarized in Figure 3.2.

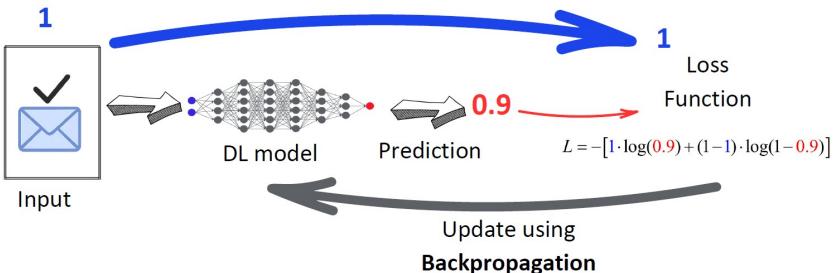


Figure 3.2: Full training flow: input → prediction → loss → backpropagation update

The model receives an input (such as an email vector), performs a forward pass to generate a prediction  $\hat{y}$ , calculates the loss  $L$ , and then updates the weights using backpropagation.

## Backpropagation

After the loss  $L$  is computed, the network must update its internal parameters-weights and biases-to reduce this loss. This is achieved using the **backpropagation** algorithm.

Backpropagation uses the **chain rule** from calculus to efficiently compute how the loss depends on each weight in the network. The chain rule allows us to break down a complex derivative into a product of simpler derivatives across layers.

*In real neural networks, there are many layers. To compute the derivative of the loss with respect to a weight in an early layer, we must chain together partial derivatives of all intermediate layers. This means calculating the derivative of the loss with respect to the output of the last layer, then that with respect to the second-to-last layer, and so on, until we reach the weight in question.*

For each parameter  $w$ , the gradient is:

$$\frac{\partial L}{\partial w}$$

This gradient tells us the sensitivity of the loss to a change in  $w$ . Intuitively:

- If  $\frac{\partial L}{\partial w} > 0$ , then increasing  $w$  increases the loss, so we should decrease  $w$
- If  $\frac{\partial L}{\partial w} < 0$ , then increasing  $w$  decreases the loss, so we may increase  $w$

## Gradient Intuition

The relationship between the loss and a weight  $w$  can be visualized as a smooth curve. At any point on this curve, the derivative (slope) indicates in which direction we should move to reduce the loss.

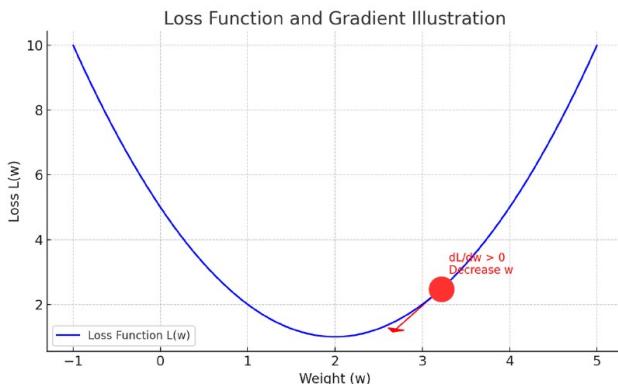


Figure 3.3: Loss as a function of weight. The red arrow shows the direction of gradient descent.

As shown in Figure 3.3, the slope of the loss curve at a given point determines how the weight should be updated.

## Gradient Descent and the Learning Rate

Once the gradients are computed, the parameters are updated using the **gradient descent** algorithm. For any parameter  $w$ , the update rule is:

$$w_{k+1} = w_k - \eta \cdot \frac{dL}{dw} \quad (3.2)$$

Where:

- $w$  is the parameter being updated at step  $k$ .
- $\eta$  is the **learning rate**, a small positive scalar
- $\frac{dL}{dw}$  is the gradient of the loss with respect to  $w$

The learning rate  $\eta$  is critical: if it is too small, training is slow. If it is too large, the model may oscillate or diverge. Choosing a good learning rate is essential to successful optimization.

*As one student put it: “I tried a learning rate of 1, and the model just bounced like a ping-pong ball.”*

## Example: One Weight Linear Model

Let us now walk through a concrete example of how gradient descent works for a very simple model with one input and one weight.

We assume a linear prediction model:

$$\hat{y} = w \cdot x \quad (3.3)$$

and a squared error loss:

$$L = (\hat{y} - y)^2 \quad (3.4)$$

Suppose we are given:

$$x = 2, \quad y = 6, \quad w = 1$$

Using Equation (3.3), the model predicts:

$$\hat{y} = 1 \cdot 2 = 2$$

Then, using Equation (3.4), we compute the loss:

$$L = (2 - 6)^2 = 16$$

Next, we compute the derivative of the loss with respect to the weight  $w$ . Taking the gradient of Equation (3.4) with respect to  $w$  using the chain rule yields:

$$\frac{dL}{dw} = 2 \cdot (wx - y) \cdot x \tag{3.5}$$

Substituting  $w = 1$ ,  $x = 2$ , and  $y = 6$  into Equation (3.5):

$$\frac{dL}{dw} = 2 \cdot (1 \cdot 2 - 6) \cdot 2 = 2 \cdot (-4) \cdot 2 = -16$$

Now we apply the gradient descent update rule from Equation (3.2):

$$w_{k+1} = w_k - \eta \cdot \frac{dL}{dw}$$

Assuming a learning rate of  $\eta = 0.1$ , we perform:

$$w_{k+1} = 1 - 0.1 \cdot (-16) = 1 + 1.6 = 2.6$$

Thus, after a single update step, the weight improves significantly - moving closer to the value that will minimize the loss.

This example demonstrates how gradient information flows from

the loss function all the way back to the weights, allowing the model to improve its predictions.

## Summary

Neural networks learn by minimizing a loss function that quantifies prediction error. Backpropagation computes how each weight contributes to the loss, using the chain rule to propagate gradients through the layers. Gradient descent then updates the weights in the direction that reduces the loss.

The learning rate controls how large each update step is. Even in deep networks, these basic principles repeat: forward pass → loss computation → backward pass → parameter update. Mastering this process is the foundation of training deep learning models.

## Lesson 4

### How Does Training Actually Work?

Up to this point, we've explored how neural networks are structured and how they generate predictions. But the core of any learning algorithm lies in its ability to improve. How exactly does a network learn from data? What happens during training, and how are the parameters - the weights and biases - updated to reduce the error?

### Batch vs. Mini-Batch Gradient Descent

In its classical form, known as **batch gradient descent**, the model processes the entire training dataset in one forward pass, computes the loss over all examples, and averages the results to form a single scalar cost:

$$L_{\text{batch}} = \frac{1}{N} \sum_{i=1}^N \ell(y_i, \hat{y}_i) \quad (4.1)$$

Here,  $N$  is the total number of training examples,  $y_i$  is the true label,  $\hat{y}_i$  is the model's prediction for input  $x_i$ , and  $\ell$  is the per-example loss function. The gradient of the loss with respect to the model parameters  $w$  is then computed:

$$\nabla_w L_{\text{batch}} = \frac{1}{N} \sum_{i=1}^N \nabla_w \ell(y_i, \hat{y}_i) \quad (4.2)$$

Finally, the weights are updated using a single gradient step:

$$w_{k+1} = w_k - \eta \cdot \nabla_w L_{\text{batch}} \quad (4.3)$$

Although this method gives a stable and precise gradient estimate, it is computationally inefficient and memory-intensive for large datasets. For this reason, batch gradient descent is rarely used in practice.

Instead, we typically use **mini-batch** gradient descent, where the dataset is divided into smaller groups of size  $B$ , called **mini-batches**. For each mini-batch, the model performs the following steps:

1. Forward passes on all  $B$  examples to compute predictions  $\hat{y}_1, \dots, \hat{y}_B$
2. Computation of the individual losses  $\ell(y_j, \hat{y}_j)$
3. Averaging the losses to compute the batch loss:

$$L_{\text{mini}} = \frac{1}{B} \sum_{j=1}^B \ell(y_j, \hat{y}_j) \quad (4.4)$$

4. Computation of the gradient of the mini-batch loss:

$$\nabla_w L_{\text{mini}} = \frac{1}{B} \sum_{j=1}^B \nabla_w \ell(y_j, \hat{y}_j) \quad (4.5)$$

5. Weight update:

$$w_{k+1} = w_k - \eta \cdot \nabla_w L_{\text{mini}} \quad (4.6)$$

*Importantly, the weights are updated **once per mini-batch**, only after averaging the loss and computing the corresponding gradient.* This approach allows for efficient use of computational resources and provides a smoother optimization trajectory compared to updating on each individual example.

## Epochs, Iterations, and Updates

To understand how training progresses over time, we introduce several standard terms:

- **Epoch** - One full pass over the entire training dataset.
- **Mini-batch** - A subset of the training data (of size  $B$ ) used to compute a single weights update.
- **Iteration** - One weights update step, performed after processing a single mini-batch.

The number of iterations per epoch is determined by:

$$\text{Iterations per epoch} = \left\lceil \frac{N}{B} \right\rceil \quad (4.7)$$

**Example:** For  $N = 10,000$  training samples and a batch size of  $B = 32$ , the number of iterations per epoch is:

$$\left\lceil \frac{10,000}{32} \right\rceil = 313 \quad (4.8)$$

This means the model performs 313 forward and backward passes per epoch, each followed by a weight update. Neural networks are typically trained over many epochs - sometimes dozens or hundreds - to allow for convergence.

*If the loss function does not decrease after several epochs, it may indicate that learning has plateaued, and training should be stopped or parameters reassessed.*

## Learning Rate

As we saw earlier, while the gradient provides the direction in which the model should change its weights, the learning rate  $\eta$  determines the step size taken in that direction. If  $\eta$  is too small, convergence is slow. If it is too large, the training process may become unstable, oscillate, or even diverge.

*The gradient points the way - but the learning rate determines how far we walk.*

Choosing an appropriate learning rate is critical to achieving efficient and stable training.

## Optimizers

The **optimizer** is the algorithm responsible for applying gradient updates to the model's parameters. It uses the gradients obtained via backpropagation, scales them by the learning rate  $\eta$ , and performs a weight update. The most basic form is **gradient descent**, introduced earlier:

$$w_{k+1} = w_k - \eta \cdot \nabla_w L_k \quad (4.9)$$

In practice, modern training workflows typically employ more advanced optimizers such as **Adam**, **RMSProp**, or **Adagrad**, which introduce mechanisms for stabilizing and accelerating convergence.

A dedicated lesson will cover the mathematical structure and trade-offs of these optimizers in detail.

## Training Dynamics Visualization

To observe how training hyperparameters influence convergence, we simulate model training on a synthetic dataset while varying the optimizer, learning rate, and batch size.

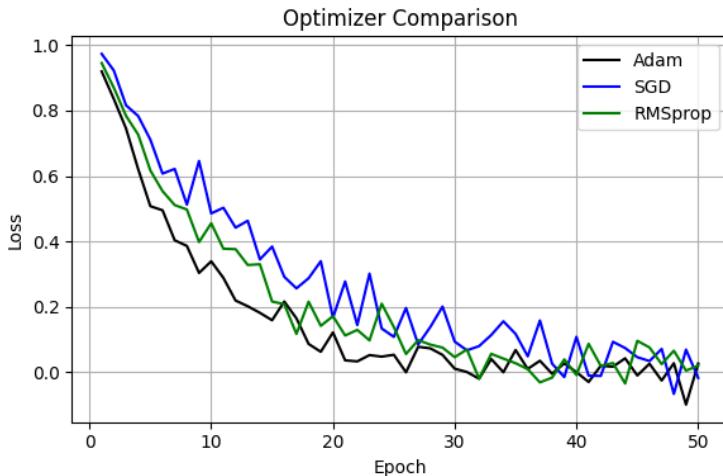


Figure 4.1: Optimizer comparison - Loss vs. Epoch. Black: Adam, Blue: SGD, Green: RMSProp.

**Optimizer Behavior:** Adam (black) achieves smooth and rapid convergence. SGD (blue) shows instability. RMSProp (green) is more stable than SGD but less efficient than Adam.

**Learning Rate Behavior:** A high learning rate (black) causes noisy or divergent updates. A low rate (green) slows down convergence. A medium value (blue) provides balance between stability and speed.

**Batch Size Behavior:** Small batches (black) produce rapid but

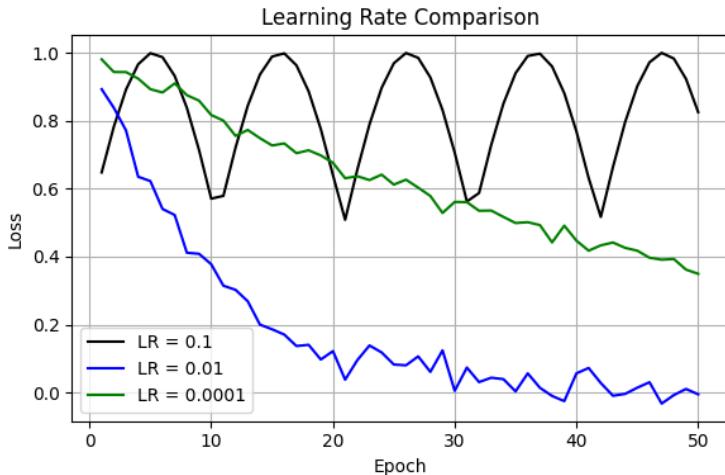


Figure 4.2: Learning rate comparison. Black: high  $\eta$ , Blue: medium  $\eta$ , Green: low  $\eta$ .

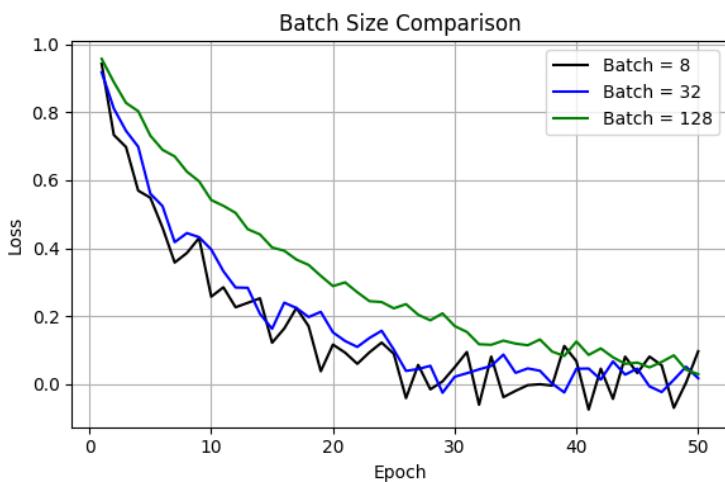


Figure 4.3: Batch size comparison. Black: small, Blue: medium, Green: large.

noisy updates. Medium sizes (blue) yield a balanced training curve. Large batches (green) result in slower but smoother convergence.

## Summary

Training dynamics are shaped by the interaction of several key components. The **batch size** determines how many examples are processed together to compute each gradient update, influencing both the stability and variance of the optimization path. The **number of epochs** defines how many full passes the model makes over the entire training dataset, affecting the overall exposure of the model to the data. The **learning rate**  $\eta$  sets the size of each update step in the parameter space, balancing convergence speed against the risk of overshooting. Finally, the **optimizer** governs the actual update rule, potentially incorporating mechanisms such as momentum or adaptive scaling to improve learning efficiency.

Effective training depends on selecting appropriate values for each of these components. Even small misconfigurations can result in unstable updates, slow convergence, or failure to learn altogether.

In the next lesson, we move from optimization to evaluation, examining how to determine whether a model has truly learned and is capable of generalizing to unseen data.

## Lesson 5

# Performance Evaluation Metrics

## Introduction

In the previous lessons, we built a neural network, trained it on real data, updated weights using gradients, and selected an optimizer. Now it is time to ask a critical question: how do we know if our model is actually good?

The goal of machine learning is not to memorize training data, but to generalize - that is, to perform well on new, unseen data. This ability is known as the model's **generalization ability**, and it is the central measure of any predictive system.

## Dataset Splitting

Before any training begins - even before choosing a model architecture - the dataset must be divided into distinct subsets. This separation is essential for ensuring that the model's performance reflects true learning rather than memorization. The **training set** is used to fit the model: gradients are computed on this data, and the parameters are updated accordingly. The **validation set** is held out during training and used to monitor performance, tune hyperparameters, and detect

signs of overfitting. Finally, the **test set** is reserved strictly for the final evaluation of generalization - it should never be used during model selection or tuning.

A common splitting strategy allocates approximately 70% of the data for training, 15% for validation, and 15% for testing. However, these ratios may vary depending on the size and nature of the dataset. The key principle is that the test set must remain untouched until the very end, serving as an unbiased estimate of real-world performance.

## Tracking Loss and Accuracy

In classification tasks, we monitor two key quantities throughout training: the **loss** and the **accuracy**. These are computed separately on the training and validation sets.

The **loss function** provides a continuous, differentiable signal indicating how far the model's predictions are from the true labels. It is sensitive to the confidence of predictions and penalizes large errors more heavily. The **accuracy**, in contrast, is a discrete metric that measures the proportion of correct predictions:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \quad (5.1)$$

Note that a lower loss does not necessarily imply higher accuracy - loss captures prediction confidence, while accuracy only reflects whether the predicted label is correct.

## Confusion Matrix and Medical Example

To better understand model behavior in binary classification, especially in sensitive domains such as healthcare, we use a tool called the

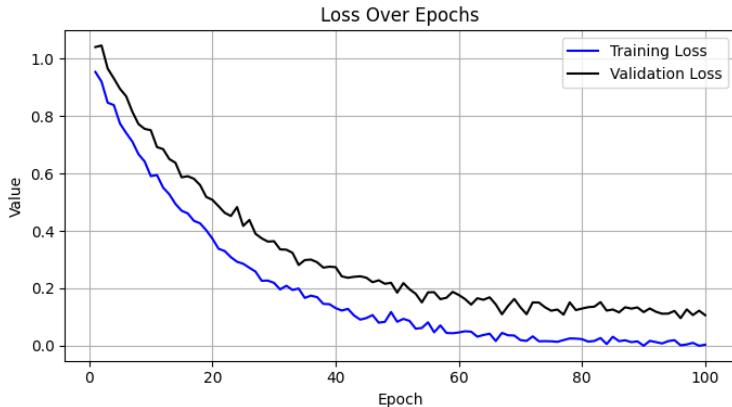


Figure 5.1: Training and validation loss over epochs: blue for training, black for validation.

**confusion matrix.** It compares predicted labels with true labels and categorizes outcomes into four cases:

### Confusion Matrix

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP)	False Negative (FN)
Actual Negative	False Positive (FP)	True Negative (TN)

The Accuracy in terms of confusion matrix components:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (5.2)$$

### Example: Disease Detection

Imagine a binary classifier built to detect a rare disease where only 1 in 100 people is sick. If the model predicts “healthy” for everyone, it achieves 99% accuracy - which is numerically high but practically disastrous. The confusion matrix helps expose such misleading cases.

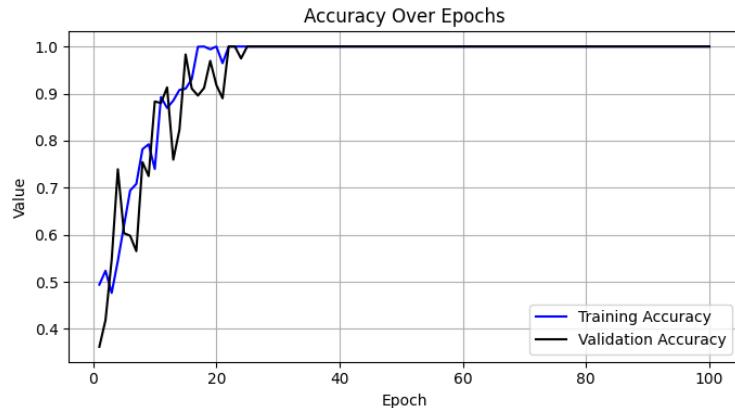


Figure 5.2: Training and validation accuracy over epochs: blue for training, black for validation.

	Predicted: Sick	Predicted: Healthy
Actually Sick	TP (Correct)	FN (Missed)
Actually Healthy	FP (False Alarm)	TN (Correct)

Each quadrant of the matrix has real-world consequences: **TP**: The disease is correctly detected. **FN**: A sick patient is missed - a critical failure. **FP**: A healthy patient is incorrectly flagged - causing unnecessary anxiety and cost. **TN**: A correct dismissal of a healthy case.

In such settings, the “positive” class refers to the condition we aim to detect - in this case, sickness. All subsequent metrics (precision, recall, F1) will be calculated with respect to the positive class.

## Precision, Recall, and F1 Score

From the confusion matrix, we define more informative classification metrics:

**Precision** - what proportion of predicted positives are correct?

$$\text{Precision} = \frac{TP}{TP + FP} \quad (5.3)$$

**Recall (Sensitivity)** - what proportion of actual positives are detected?

$$\text{Recall} = \frac{TP}{TP + FN} \quad (5.4)$$

**F1 Score** - harmonic mean of precision and recall:

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5.5)$$

These metrics are especially important in imbalanced datasets. High precision with low recall indicates a conservative model; high recall with low precision implies many false alarms.

## ROC Curve and AUC

When models output probabilities, we apply a decision threshold (typically 0.5) to produce a binary classification. Adjusting this threshold affects sensitivity and specificity.

The **ROC curve** plots the True Positive Rate (Recall) against the False Positive Rate for various thresholds. The area under this curve - the **AUC** - summarizes the model's ability to distinguish between classes:

- $\text{AUC} = 1.0 \rightarrow$  perfect classification
- $\text{AUC} = 0.5 \rightarrow$  no better than random guessing
- $\text{AUC} < 0.5 \rightarrow$  inverted predictions

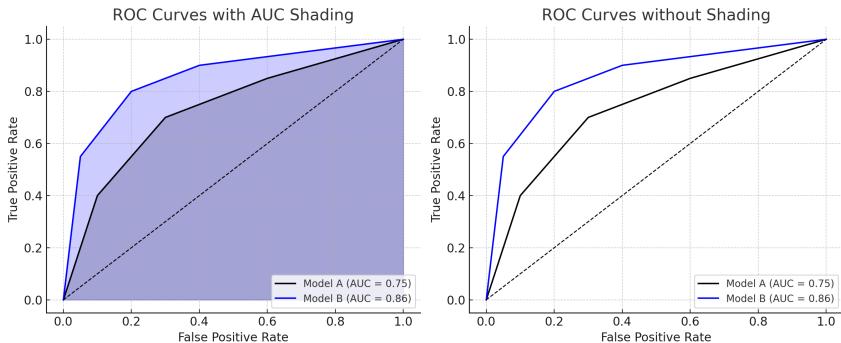


Figure 5.3: ROC curve with AUC (shading): Model A: black, AUC=0.75; Model B: blue, AUC=0.86.

## Regression Metrics

When the model outputs a continuous value rather than a class, we evaluate it using **regression metrics**. These quantify the discrepancy between predicted values  $\hat{y}_i$  and actual values  $y_i$ .

### Mean Absolute Error (MAE)

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (5.6)$$

MAE measures the average magnitude of errors, treating all deviations equally. It is robust to outliers and shares the same units as the output variable.

### Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (5.7)$$

MSE emphasizes larger errors due to squaring. It is commonly used as a loss function due to its smooth, differentiable nature, but is more

sensitive to outliers.

## Coefficient of Determination ( $R^2$ )

The **coefficient of determination**, denoted by  $R^2$ , is a widely used metric for evaluating the performance of regression models. It measures how well the predicted values  $\hat{y}_i$  approximate the actual label values  $y_i$ , compared to a simple baseline: always predicting the average of the observed outcomes.

Formally,  $R^2$  is defined as:

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} \quad (5.8)$$

Here,  $SS_{\text{res}}$  is the residual sum of squares - the total squared error between the predicted values and the true values.  $SS_{\text{tot}}$  is the total sum of squares - the total variance in the true values relative to their mean. Specifically,  $\bar{y}$  denotes the mean of the true target values across the dataset, that is:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

This mean serves as a naïve baseline predictor.

$$SS_{\text{res}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad SS_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

The value of  $R^2$  reflects how much of the original variance in the data has been explained by the model. A value of  $R^2 = 1$  indicates a perfect fit: the model predictions exactly match the true values. A value of  $R^2 = 0$  means the model performs no better than simply predicting the mean  $\bar{y}$  for every instance. If  $R^2 < 0$ , the model is worse than the baseline and introduces additional error.

## Additional Metrics for Classification

### Balanced Accuracy

Balanced Accuracy accounts for imbalanced datasets by averaging recall across classes:

$$\text{Balanced Accuracy} = \frac{1}{2} \left( \frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right) \quad (5.9)$$

This ensures that both classes contribute equally to the final score, regardless of prevalence.

### Matthews Correlation Coefficient (MCC)

$$\text{MCC} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (5.10)$$

MCC is a balanced measure even for imbalanced data. It considers all four outcomes and behaves like a correlation coefficient:  $+1$  indicates perfect prediction,  $0$  indicates random guessing, and  $-1$  means total disagreement.

## Summary

Evaluating model performance requires more than a single number. In classification, accuracy is a starting point, but deeper insight comes from precision, recall, F1 score, and AUC. In regression, metrics like MAE, MSE, and  $R^2$  quantify the quality of numerical predictions. Metrics such as balanced accuracy and MCC provide robustness in imbalanced settings. Above all, evaluation should be performed on the held-out test set to ensure a reliable estimate of generalization performance.

## Lesson 6

### Overfitting and Regularization

#### Introduction

In the early stages of training a neural network, everything may appear to be going well: the training loss decreases steadily, and the training accuracy improves epoch by epoch. However, when new, unseen data is introduced, the model's performance can suddenly collapse. This phenomenon is not a bug - it is a fundamental challenge in deep learning, known as **overfitting**.

Overfitting occurs when the model fits the training data too closely, capturing not only the underlying patterns but also noise, anomalies, and incidental correlations. As a result, the model performs well on the data it has seen, but fails to generalize to new examples. In contrast, **underfitting** describes a scenario where the model fails to capture even the patterns in the training data, often due to insufficient capacity or inadequate training.

Both overfitting and underfitting degrade the model's **generalization ability** - the fundamental goal.

## Conceptual Example

Consider a student who memorizes answers from last year's exam without understanding the concepts. When presented with a new version of the exam with slightly altered questions, the student is unable to answer correctly. The analogy is clear: memorization is not learning.

In a similar way, a neural network may continue to decrease its training loss while the validation loss begins to rise - a classic indicator of overfitting. Initially, the model improves on both training and validation sets, but at a certain point, it begins to *memorize* rather than learn.

## Loss Behavior as a Diagnostic Tool

During training, we monitor the loss on both the training and validation sets. Let  $L_{\text{train}}(t)$  and  $L_{\text{val}}(t)$  denote the training and validation loss, respectively, at epoch  $t$ .

In the early epochs:

$$L_{\text{train}}(t) \downarrow, \quad L_{\text{val}}(t) \downarrow$$

At the point of overfitting onset:

$$L_{\text{train}}(t) \downarrow, \quad L_{\text{val}}(t) \uparrow$$

This divergence is a critical signal that the model is no longer learning meaningful structure, but rather memorizing the training data.

# Why Overfitting Happens

Modern neural networks often contain millions of parameters. With enough capacity, a model can always achieve near-zero training loss - even on randomly labeled data. However, such a solution captures the data's noise rather than its underlying structure.

Perfect accuracy on the training set is often a red flag, not a cause for celebration. It usually indicates over-parametrization and a lack of inductive bias.

## Three Core Techniques for Preventing Overfitting

### 1. Dropout

**Dropout** is a simple yet powerful regularization technique. During training, each hidden unit is randomly *dropped* - i.e., temporarily deactivated - with probability  $p$ . As a result, different subnetworks are trained in each forward pass.

Formally, let  $h^{(l)}$  be the activations of the  $l$ -th layer. Dropout introduces a binary mask  $m^{(l)} \sim \text{Bernoulli}(1 - p)$ :  $\tilde{h}^{(l)} = m^{(l)} \cdot h^{(l)}$ . At inference time, all neurons are active, and weights are scaled appropriately to account for the absence of dropout.

Dropout prevents co-adaptation of neurons and encourages the network to develop redundant, robust representations. The result is improved generalization.

### 2. Early Stopping

**Early stopping** monitors the validation loss during training. If no improvement is observed for a fixed number of epochs (called *patience*),

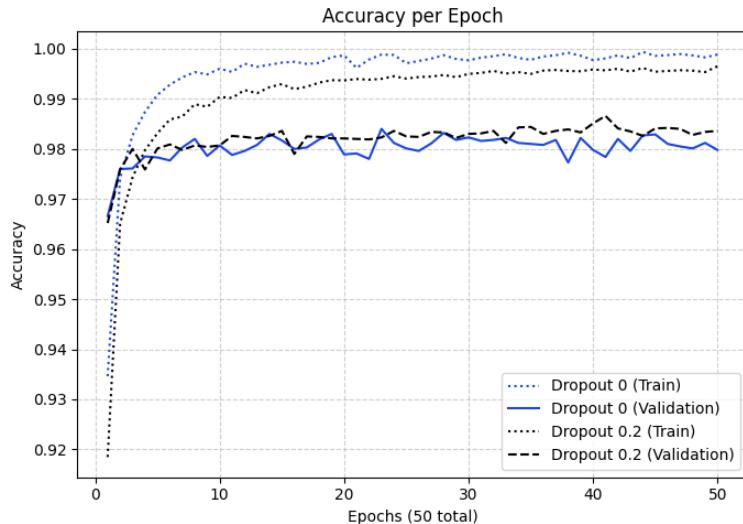


Figure 6.1: Accuracy over epochs. Blue: without dropout. Black dashed: dropout with  $p = 0.2$ .

training is terminated early.

This simple strategy prevents the model from continuing to optimize the training loss at the expense of generalization. It is especially useful when the ideal number of training epochs is unknown.

### 3. Weight Decay (L2 Regularization)

**L2 regularization**, also known as *weight decay*, penalizes large weights by adding a quadratic term to the loss function:

$$L_{\text{reg}} = L + \lambda \sum_i w_i^2 \quad (6.1)$$

Here,  $L$  is the original loss (e.g., cross-entropy),  $w_i$  denotes the  $i$ -th weight, and  $\lambda$  is a regularization coefficient.

The intuition is clear: large weights imply a strong dependency on

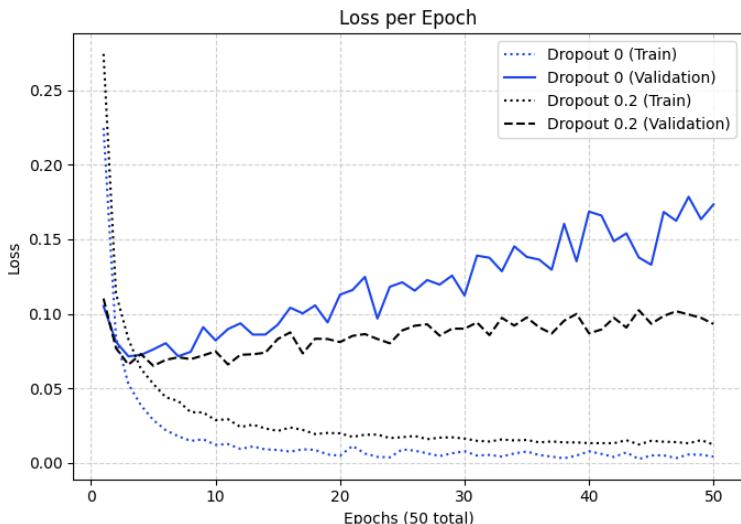


Figure 6.2: Validation loss and over epochs. Blue: without dropout.  
Black dashed: dropout with  $p = 0.2$ .

specific features. Penalizing such weights encourages the network to distribute importance more evenly, thereby promoting more stable and generalizable learning.

## When to Apply Regularization?

In many real-world scenarios, especially when the validation set is small or noisy, it may be difficult to detect overfitting with certainty. For this reason, it is often recommended to apply regularization techniques by default - as a preventive measure, not merely as a response to observed failure.

Technique	Mechanism	Purpose
Dropout	Randomly deactivates neurons during training	Prevents co-adaptation, encourages redundancy
Early Stopping	Halts training when validation loss stops improving	Avoids over-optimization beyond generalization point
L2 Regularization (Weight Decay)	Adds penalty term $\lambda \sum w_i^2$ to the loss function	Discourages large weights, encourages smooth solutions

Table 6.1: Comparison of common regularization techniques.

## Summary

In this lesson, we examined one of the most important challenges in deep learning: overfitting. We defined it, diagnosed it using loss curves, and introduced three powerful tools to combat it: dropout, early stopping, and weight decay.

These techniques serve not just to constrain models, but to guide them toward learning meaningful patterns rather than memorizing noise. In the next module, we will explore advanced training techniques and architectures that allow for deeper, more expressive networks without sacrificing generalization.

## Lesson 7

### Why Do We Need Convolution?

#### Introduction

In previous lessons, we examined how neural networks learn from structured data. In this lesson, we focus on images - a domain with strong spatial structure. Fully connected networks (MLPs) ignore this structure and require an enormous number of parameters (weights). Convolutional Neural Networks (CNNs) offer a more efficient and scalable solution. But why is convolution necessary at all?

#### Fully Connected Layers and Their Limitations

Let us consider a small RGB image of size  $32 \times 32$  pixels. The number of input features is:

$$3 \times 32 \times 32 = 3,072 \quad (7.1)$$

If this flattened vector is connected to a dense layer of 100 neurons, the number of weights in just the first layer is:

$$W_{Total} = 3072 \times 100 = 307,200 \quad (7.2)$$

This dense connection loses spatial awareness - the network does

not know that pixels next to each other are related. As we add more layers, the parameter count grows rapidly and might lead to overfitting.

**Table 7.1: Parameter count and spatial awareness**

Architecture	Parameters	Spatial Awareness
MLP (100 neurons)	307,200	✗
CNN (32 filters of $3 \times 3$ )	864	✓

### Why Convolution Works - Inductive Bias

There's a well-known result in machine learning theory called the *No Free Lunch* theorem. It says that if you average the performance of an algorithm over all possible data-generating functions, no method beats another. The only reason CNNs outperform MLPs on images is because we assume **locality and translation invariance**. Likewise, there are architectures (Transformers) that work well on text because of their ability to model relationships between any pair of tokens.

In practice, this means that there is no universally best model, only models that match the structure of your data. Convolutional architectures succeed in vision tasks precisely because they encode assumptions about how visual patterns behave.

## The Convolution Operator

Convolution applies a small learnable filter across local image regions. Let  $I \in \mathbb{R}^{H \times W}$  be a grayscale input image, and let  $K \in \mathbb{R}^{k \times k}$  be a kernel (filter). The 2D discrete convolution is defined as:

$$F(i, j) = \sum_{u=1}^k \sum_{v=1}^k K(u, v) \cdot I(i + u - 1, j + v - 1) \quad (7.3)$$

Where:

- $F(i, j)$  is the output feature map value at location  $(i, j)$ .
- $K(u, v)$  is the kernel weight at position  $(u, v)$ .
- $I(i + u - 1, j + v - 1)$  is the input pixel under the kernel window.

This operation is repeated across all spatial positions, reusing the same kernel - which greatly reduces the number of trainable parameters.

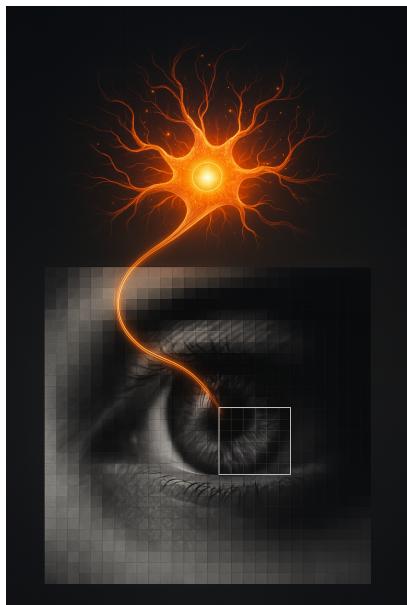


Figure 7.1: Filter applied to eye region - local pattern detection.

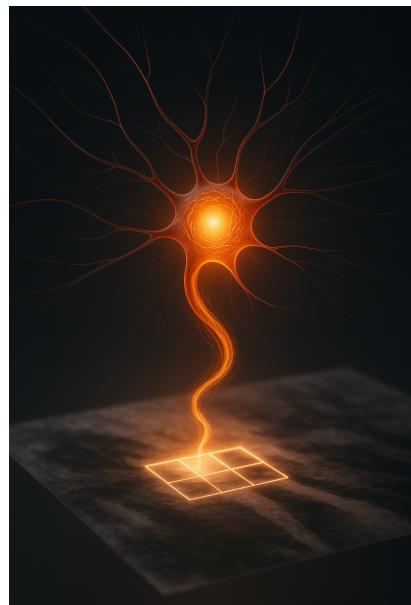


Figure 7.2: Filter reused - translation invariance.

## Convolutional Intuition

Convolution extracts reusable patterns from local regions. For example, a filter may detect "eye" or "corner" features that appear anywhere in the image. Just as the human brain doesn't memorize entire objects but rather learns their components (e.g., curves, edges, eyes), CNNs specialize in recognizing spatially invariant patterns.

## Properties

One of the core advantages of convolution is its use of **parameter sharing**, where the same set of weights is applied across different spatial locations. This greatly reduces the number of learnable parameters and enhances computational efficiency. In addition, convolutional models are naturally robust to translation - they can recognize objects regardless of where they appear in the image - leading to improved generalization and overall model stability. Although convolution can reduce spatial dimensions, its behavior is also controlled by:

- **Stride ( $s$ )**: the number of pixels the filter shifts each step.
- **Padding ( $p$ )**: how many border pixels are added to preserve output size.

These will be formally defined in the next lesson.

## The CIFAR-10 Dataset

To compare MLPs and CNNs, we use the **CIFAR-10** benchmark - a standard dataset in image classification. It contains 60,000 color images of size  $32 \times 32 \times 3$ , across 10 classes such as airplane, dog, and automobile.

CIFAR-10 is frequently used to test model architectures because:



Figure 7.3: Sample images from CIFAR-10.

- It is small enough to allow rapid experimentation.
- It includes real-world variation across 10 diverse classes.
- State-of-the-art CNNs can achieve over 90% accuracy on it.

## Case Study: MLP vs. CNN on CIFAR-10

We trained two models on CIFAR-10 for 30 epochs:

1. A fully connected MLP with 2 dense layers and ReLU activation.
2. A CNN with two convolutional layers and ReLU activation.

### Observations:

- MLP reached only  $\sim 50\%$  accuracy.
- CNN reached  $\sim 68\%$ , significantly better.

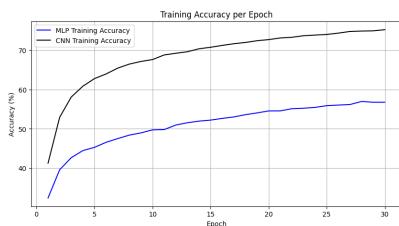


Figure 7.4: Training accuracy: CNN (black) reaches higher accuracy than MLP (blue), faster.

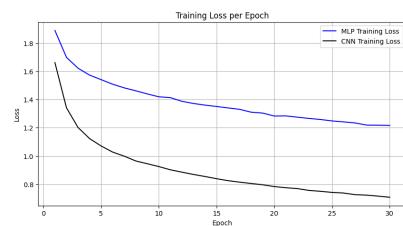


Figure 7.5: Training loss: CNN (black) reaches lower error than MLP (blue), faster.

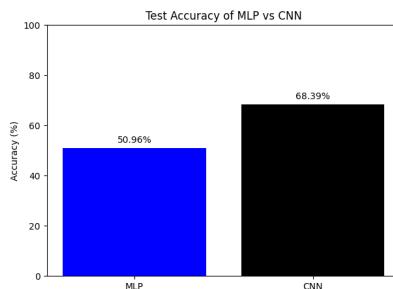


Figure 7.6: Training loss - CNN converges faster and reaches lower error.

## Summary

Convolution introduces a powerful inductive bias tailored for visual data. Unlike fully connected networks that disregard spatial structure, Convolutional Neural Networks (CNNs) preserve the locality of features by processing data in small, overlapping regions. This allows them to capture meaningful spatial patterns such as edges, textures, and shapes. Through the mechanism of weight sharing, CNNs require significantly fewer parameters, which not only improves generalization but also makes training more efficient. Another critical advantage is their ability to detect patterns regardless of where they appear in the image, offering a form of translation robustness. These properties make CNNs highly scalable and effective for real-world image data, even as input dimensions grow.

In the next lesson, we will formalize the mathematical definitions of stride, padding, and pooling, and begin constructing complete convolutional architectures for various tasks.

## Lesson 8

### How Does a CNN Work?

#### Introduction

Now that we understand why convolution is essential, we turn our focus to the internal structure of a Convolutional Neural Network (CNN). In this lesson, we define the components of a convolutional layer, explain how filters slide across the image, and introduce critical mechanisms such as stride, padding, and pooling.

CNNs operate by extracting local features using learned filters. These filters convolve across the image and produce feature maps - each emphasizing different patterns such as edges, corners, or textures. When stacked hierarchically, convolutional layers enable deep semantic understanding of visual data.

#### Convolutional Layer: Filters and Sliding Window

Each convolutional layer consists of multiple filters, also called kernels. Each filter is a learnable matrix of weights, typically of size  $k \times k$ , for instance  $3 \times 3$ .

At each spatial location, the filter performs element-wise multiplication with the corresponding input patch, sums the result, and places it into the output feature map. This operation is called a sliding window convolution.

Formally, let  $I \in \mathbb{R}^{H \times W}$  be a grayscale image and  $K \in \mathbb{R}^{k \times k}$  a convolution kernel. The output  $F \in \mathbb{R}^{H_o \times W_o}$  is defined by:

$$F(i, j) = \sum_{u=1}^k \sum_{v=1}^k K(u, v) \cdot I(i + u - 1, j + v - 1) \quad (8.1)$$

Here,  $F(i, j)$  is the output value at position  $(i, j)$ ,  $K(u, v)$  is the filter coefficient at  $(u, v)$ , and  $I(i + u - 1, j + v - 1)$  is the input pixel under the filter.

Applying  $M$  such filters produces  $M$  feature maps - one for each pattern the network has learned to detect.

## Padding and Border Effects

Without padding, the filter cannot fully cover the border regions of the image. As a result, each convolution reduces the spatial size of the feature map. To preserve dimensions, we apply **zero-padding** - adding rows and columns of zeros around the input.

$$x = \begin{bmatrix} ? & ? \\ 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 3 & 4 & 0 \\ 0 & 5 & 6 & 7 & 8 & 0 \\ 0 & 9 & 10 & 11 & 12 & 0 \\ 0 & 13 & 14 & 15 & 16 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Figure 8.1: Zero-padding: added borders allow full coverage of input dimensions.

Padding is especially important in deep networks, where multiple

convolutional layers might otherwise shrink the spatial resolution excessively.

## Stride and Its Effect

The **stride**  $s$  determines the number of pixels the filter moves at each step. A stride of 1 means the filter moves one pixel at a time, resulting in high-resolution output. A stride of 2 causes the filter to skip alternate positions, reducing output dimensions and increasing computational efficiency.

$$\begin{array}{ll} \text{Stride } = 1 & x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \\ \text{Stride } = 2 & x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \end{array}$$

Figure 8.2: Effect of stride: top – stride 1 (dense scanning), bottom – stride 2 (coarser sampling).

Now, combining the effects of input height  $H$ , padding  $p$ , kernel size  $k$ , and stride  $s$ , the output dimension is:

$$H_{\text{out}} = \left\lfloor \frac{H + 2p - k}{s} + 1 \right\rfloor \quad (8.2)$$

$$W_{\text{out}} = \left\lfloor \frac{W + 2p - k}{s} + 1 \right\rfloor \quad (8.3)$$

## Pooling: Max and Average Variants

After convolution, we often apply a pooling operation to downsample the feature map. The most common form is **max pooling**, which selects the highest value from each local patch. Alternatively, **average pooling** computes the mean of values in the region.

Formally, if  $R_{i,j}$  denotes the pooling region for output location  $(i, j)$ , then:

**Max pooling:**

$$P_{\text{max}}(i, j) = \max_{(u,v) \in R_{i,j}} F(u, v) \quad (8.4)$$

**Average pooling:**

$$P_{\text{avg}}(i, j) = \frac{1}{|R_{i,j}|} \sum_{(u,v) \in R_{i,j}} F(u, v) \quad (8.5)$$

$$x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow 4$$

$$x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow 2.5$$

Figure 8.3: Pooling variants: top - max pooling; bottom - average pooling. Both reduce spatial resolution.

## CNN Block: A Basic Pipeline

A typical CNN block consists of the following stages:

- **Input:** an image, e.g., of size  $28 \times 28$  pixels.

- **Convolution:** apply  $3 \times 3$  filters (e.g., 8 learnable filters).
- **Activation:** apply a non-linear function, typically ReLU.
- **Pooling:** reduce spatial dimensions using max pooling.
- **Flatten:** convert the resulting feature maps into a 1D vector.
- **Fully Connected Layer:** perform the final classification based on the extracted features.

These blocks are often stacked to form deeper models. With each layer, the receptive field of the output grows.

## Receptive Field and Depth

Although each convolutional unit sees only a small region, stacking multiple layers increases the **receptive field** - the region of the input that influences a neuron.

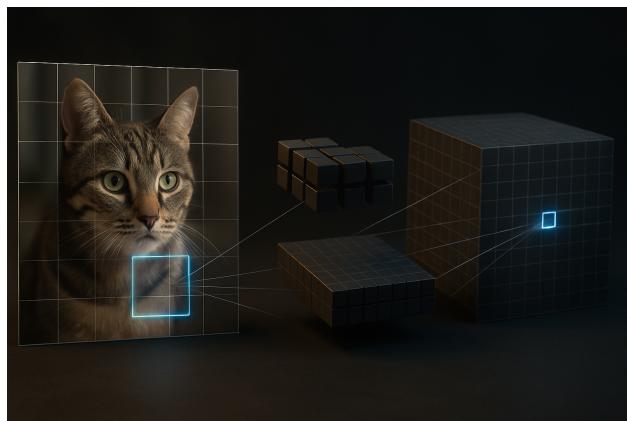


Figure 8.4: Receptive field grows with depth - allowing the model to capture global context.

In a well-designed CNN, lower layers detect local edges, mid-level layers identify textures and motifs, and upper layers respond to full objects.

## Summary

A CNN combines filters, stride, padding, and pooling into a powerful hierarchical system for learning spatial features. These mechanisms allow the network to reduce data size, preserve structure, and build progressively more abstract representations of visual input. In the next lesson, we introduce architectural enhancements like batch normalization and dropout that further improve performance and robustness.

## Lesson 9

# Sequences and Time: RNN, GRU, and LSTM

## Introduction

Up until now, we dealt with inputs such as images, where all data was available at once. However, many domains—such as natural language, music, and sensor streams—require processing **sequential data**. In such tasks, temporal context is essential: the model must not only “see” but also **remember**.

In this lesson, we introduce the family of Recurrent Neural Networks (RNNs), starting with the basic RNN, analyzing its limitations, and then presenting two widely used architectures: GRU and LSTM. We conclude with a powerful enhancement—bi-directional LSTM.

## Recurrent Neural Networks (RNN)

Unlike images, where all pixels are processed simultaneously, sequences unfold step by step. At each time step  $t$ , a new input  $x_t$  arrives, and the model must update its internal state to reflect both the current input and past history.

An RNN maintains a hidden state  $h_t$  that is updated recurrently as follows:

$$h_t = f_h(\mathbf{V}x_t + \mathbf{W}h_{t-1} + b_h) \quad (9.1)$$

$$\hat{y}_t = f_y(\mathbf{U}h_t + b_y) \quad (9.2)$$

Where:

- $x_t$  is the input at time  $t$ .
- $h_t$  is the hidden state at time  $t$ .
- $\mathbf{V}$ ,  $\mathbf{W}$ , and  $\mathbf{U}$  are learnable weight matrices.
- $b_h$ ,  $b_y$  are biases.
- $f_h$  is the activation function (e.g., tanh or ReLU).
- $\hat{y}_t$  is the model's output at time  $t$ .

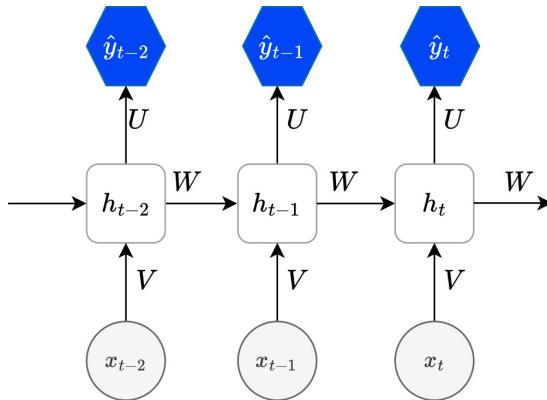


Figure 9.1: Unrolled Recurrent Neural Network architecture.

This mechanism enables information to flow through time. However, classical RNNs suffer from the **vanishing gradient problem**, especially over long sequences. Gradients diminish during backpropagation through time (BPTT), limiting the network's ability to learn long-term dependencies.

## LSTM: Long Short-Term Memory

The Long Short-Term Memory network (LSTM) introduces a separate internal memory cell  $C_t$  and a set of **gating mechanisms** that regulate the flow of information. This architecture was designed to mitigate the vanishing gradient problem. LSTM updates are defined as follows:

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i) \quad (\text{input gate}) \quad (9.3)$$

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f) \quad (\text{forget gate}) \quad (9.4)$$

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \quad (\text{output gate}) \quad (9.5)$$

$$\tilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c) \quad (9.6)$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad (9.7)$$

$$h_t = o_t \odot \tanh(C_t) \quad (9.8)$$

Here:  $\sigma$  is the sigmoid activation function,  $\odot$  denotes element-wise multiplication,  $C_t$  is the cell state, and  $h_t$  is the hidden state.

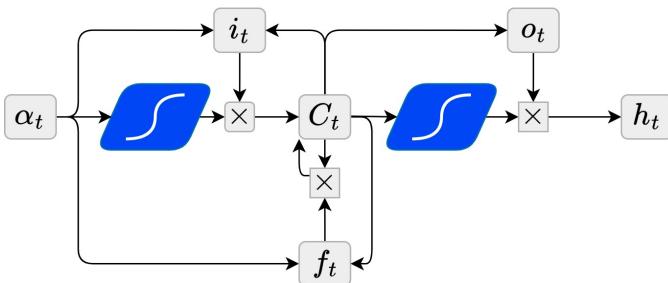


Figure 9.2: Internal structure of an LSTM unit.

This design enables selective memorization and forgetting, making LSTM ideal for capturing both short and long-term dependencies.

## GRU: Gated Recurrent Unit

GRU is a simplified variant of LSTM that merges the input and forget gates into a single update gate. It has fewer parameters and often trains faster while maintaining similar performance. The GRU equations are:

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \quad (\text{update gate}) \quad (9.9)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad (\text{reset gate}) \quad (9.10)$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h) \quad (9.11)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \quad (9.12)$$

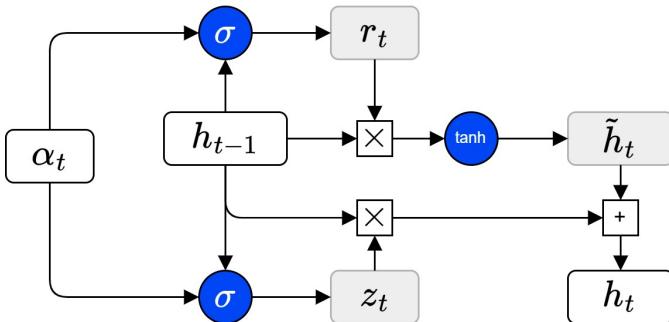


Figure 9.3: GRU architecture showing two gates: update and reset.

GRUs are computationally cheaper and are often preferred for long sequences when training time is a concern.

## Bi-Directional LSTM

Some tasks benefit from having access to both past and future context. A Bi-directional LSTM (Bi-LSTM) achieves this by running two

LSTMs in parallel: One processes the sequence forward (from  $t = 1$  to  $t = T$ ), and one processes it backward (from  $t = T$  to  $t = 1$ ). Their outputs are concatenated at each time step:

$$h_t^{\text{bi}} = [h_t^{\rightarrow}; h_t^{\leftarrow}] \quad (9.13)$$

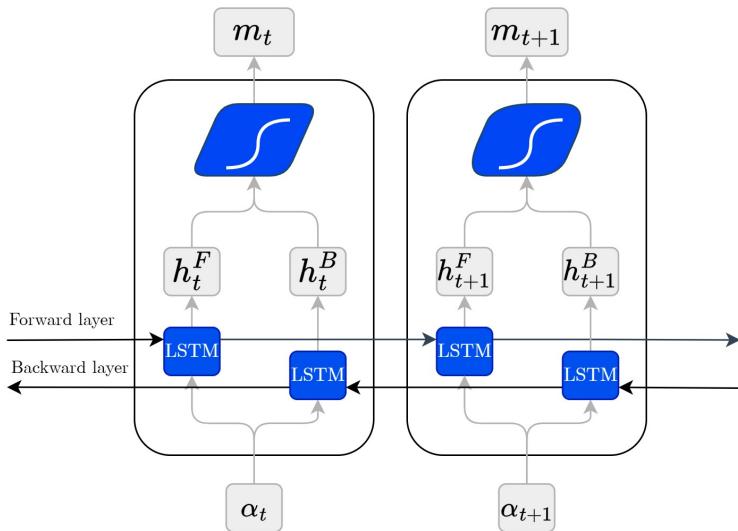


Figure 9.4: Bi-directional LSTM: combining forward and backward hidden states.

Bi-LSTMs are particularly effective in text processing tasks, where the meaning of a word often depends on both previous and following words.

## Example: Sentiment Analysis

Consider the sentence:

*“The service was okay, but the food was absolutely terrible.”*

A left-to-right RNN might process the early part — “*The service was okay*” — and conclude that the sentiment is neutral or mildly positive. However, the key sentiment signal arrives only at the end: “*absolutely terrible*”, indicating strong negativity.

A unidirectional RNN may fail to revise its earlier prediction. In contrast, a **Bidirectional LSTM** processes the sequence in both directions. It can incorporate context from the end (“*terrible*”) into its understanding of the whole sentence, leading to a more accurate sentiment classification. This demonstrates the network’s reliance on context.

Model	Gates	Memory	Pros	Use
RNN	None	Short-term	Simple	Short seq.
LSTM	In, Forget, Out	Long-term	Stable, powerful	Text, audio
GRU	Update, Reset	Long-term	Fast, fewer params	Same as LSTM

Table 9.1: Summary of recurrent architectures.

## Summary

Recurrent models are essential for tasks involving temporal data. The basic RNN introduces the principle of memory across time. LSTM and GRU improve upon this by enabling the network to retain and control information more effectively.

For tasks requiring future context, Bi-LSTM extends this capability. However, these architectures still rely on sequential computation, which can be a limitation for scalability.

In the next lessons, we will introduce the **Transformer** architecture - a revolutionary framework that replaces recurrence with self-attention and dominates the state-of-the-art in language modeling and sequence tasks.

## Lesson 10

### Autoencoders for Dimensionality Reduction

#### Introduction

Until now, we have focused on supervised learning - training models on labeled input-output pairs. In this lesson, we explore a different paradigm: **unsupervised learning**, where no labels are given, and the goal is to uncover hidden structure within the data.

Machine learning is typically categorized into three main types: (1) *Supervised learning*, where the model maps inputs  $x$  to known outputs  $y$ ; (2) *Unsupervised learning*, where the task is to find patterns or compress data without labels; (3) *Reinforcement learning*, where the agent learns through interaction and reward.

While this course focuses on supervised learning, unsupervised models are equally essential. Among them, the **Autoencoder (AE)** plays a central role in modern deep learning. An autoencoder is a neural network trained to reconstruct its own input - and through this process, it learns an internal compressed representation of the data.

#### Autoencoder Architecture

An autoencoder is composed of two sub-networks:

- **Encoder**  $f_{\text{enc}}$ : compresses the input  $x \in \mathbb{R}^d$  into a latent vector  $z \in \mathbb{R}^k$ , where  $k \ll d$ .
- **Decoder**  $f_{\text{dec}}$ : reconstructs the input  $\hat{x} \in \mathbb{R}^d$  from  $z$ .

Formally, the encoder and decoder satisfy:

$$z = f_{\text{enc}}(x) \quad (10.1)$$

$$\hat{x} = f_{\text{dec}}(z) \quad (10.2)$$

The network is trained to minimize the reconstruction error:

$$\mathcal{L}(x, \hat{x}) = \|x - \hat{x}\|^2 = \sum_{i=1}^d (x_i - \hat{x}_i)^2 \quad (10.3)$$

This loss encourages the model to retain essential information in  $z$ , while discarding redundancy.

## Learning Compact Representations

By enforcing a bottleneck (i.e.,  $k \ll d$ ), the autoencoder is constrained to learn only the most meaningful patterns in the input. The encoder extracts abstract features, and the decoder learns to reassemble them into an accurate approximation of the original input. It is particularly powerful for high-dimensional data such as images, audio, and sensor signals, where interpretability and compression are valuable.

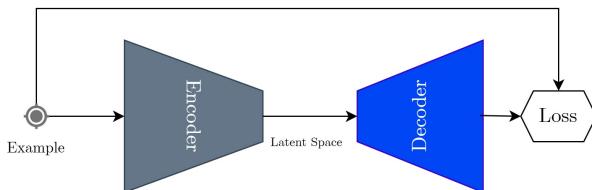


Figure 10.1: Autoencoder Illustration.

## Example: Dimensionality Reduction on MNIST

To demonstrate this idea, we apply an autoencoder to the MNIST dataset - consisting of  $28 \times 28 = 784$  pixel grayscale images of handwritten digits. We design an autoencoder with:

- Input:  $x \in \mathbb{R}^{784}$
- Latent space:  $z \in \mathbb{R}^2$
- Decoder:  $\hat{x} \in \mathbb{R}^{784}$

The total number of trainable parameters in this network is 218,514.

Layer (type)	Output Shape	Param #
input_layer (InputLayer)	(None, 784)	0
dense (Dense)	(None, 128)	100,480
dense_1 (Dense)	(None, 64)	8,256
dense_2 (Dense)	(None, 2)	130
dense_3 (Dense)	(None, 64)	192
dense_4 (Dense)	(None, 128)	8,320
dense_5 (Dense)	(None, 784)	101,136

Figure 10.2: Architecture of the MNIST autoencoder. Total trainable parameters: 218,514.

After training, we extract the latent vectors  $z$  for all digits. Even though no labels were provided during training, we find that digits with similar structure (e.g., "1", "7") cluster together.

This demonstrates that the model has captured the internal structure of the data - a hallmark of successful representation learning.

## Reconstruction Results

We visualize several digit reconstructions from the autoencoder. For each input image  $x$ , we pass it through the encoder-decoder pipeline and compare the reconstruction  $\hat{x}$ :

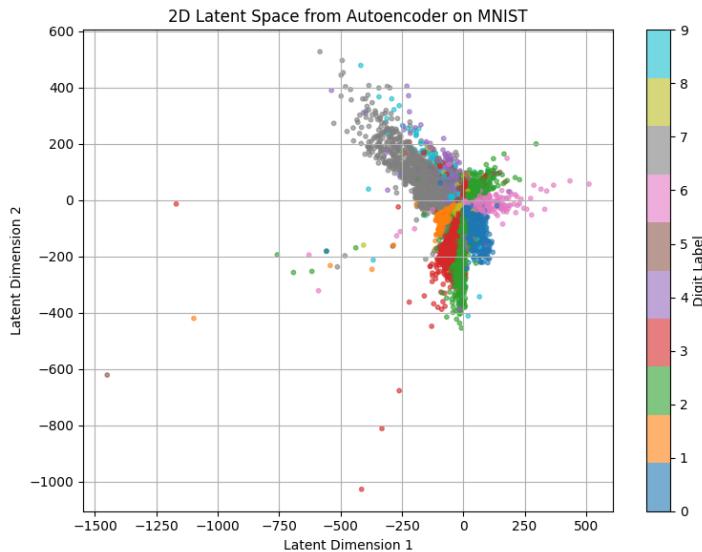


Figure 10.3: 2D latent space learned by the autoencoder. Colors indicate true digit labels (used for visualization only).

Despite the bottleneck, the model manages to preserve most of the perceptual structure, showing it has extracted meaningful latent features.

## Training Curve

The following figure shows the reconstruction loss (mean squared error) as a function of training epochs. As training progresses, the loss decreases steadily:

## Comparison with PCA

Principal Component Analysis (PCA) is a linear technique for dimensionality reduction. It finds orthogonal directions that maximize variance. However, PCA cannot model nonlinearities in the data.



Figure 10.4: Original (top) and reconstructed (bottom) digits using a trained autoencoder for 300 epochs.

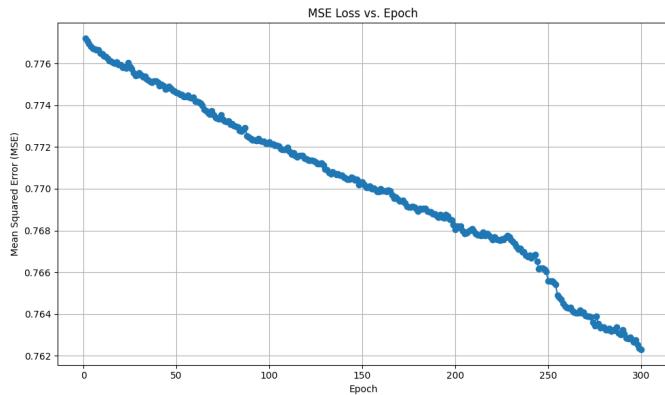


Figure 10.5: MSE reconstruction loss vs. epochs. The model converges stably to a low error.

Autoencoders generalize PCA using nonlinear functions (e.g., ReLU activations), enabling them to learn curved manifolds. The key difference is that PCA finds a linear subspace minimizing squared reconstruction error while autoencoders learn a nonlinear function minimizing the same loss. Autoencoders are more flexible and better suited for structured data such as images.

## Applications and Extensions

Autoencoders have been successfully employed in a wide range of deep learning tasks. Some of their primary applications include:

- Learning low-dimensional embeddings from high-dimensional inputs.
- Signal compression for efficient storage or transmission.
- Anomaly detection based on elevated reconstruction error.
- Pretraining of deep neural networks, followed by supervised fine-tuning.
- Denoising of inputs by removing structured or stochastic noise.

Beyond the standard autoencoder, several extensions have been developed to address specific modeling needs. One of them is the Variational Autoencoder (VAE), where the encoder outputs the parameters of a Gaussian distribution - a mean vector  $\mu(x)$  and a standard deviation vector  $\sigma(x)$  - from which the latent variable  $z$  is sampled as:

$$z \sim \mathcal{N}(\mu(x), \sigma^2(x)). \quad (10.4)$$

## Summary

Autoencoders are a fundamental unsupervised learning tool that reveals the internal structure of data by compressing and reconstructing it. They do not rely on labels or external supervision, yet they learn meaningful representations that can be visualized, clustered, and reused in downstream tasks.

The MNIST example illustrates how autoencoders uncover geometric patterns in digit space without ever being told what a digit is. Their strength lies in the ability to generalize beyond linear projections and learn the essence of data structure. This makes them indispensable for many modern deep learning pipelines.

## Lesson 11

### Self-Attention and the Transformer Principle

#### Introduction

The Transformer architecture, introduced in 2017 in the paper *Attention Is All You Need*, revolutionized deep learning for sequential data. Originally developed for natural language processing (NLP), the Transformer quickly became the foundation of modern generative models such as GPT and BERT. Its core innovation - self-attention - replaced recurrence with a mechanism that allows every token to attend to every other, enabling greater parallelism and better modeling of long-range dependencies. Later, the same principle was extended to other domains: Vision Transformers (ViT), for example, apply attention over image patches.

#### Why RNNs Are Not Enough

Recurrent Neural Networks (RNNs) and their variants like LSTMs process sequences token by token, updating a hidden state that stores past information. This inherently sequential approach makes parallelization difficult and limits their ability to capture distant dependencies.

**Example:** “*The book that the professor recommended at*

*the end of the lecture was truly brilliant.”* To correctly interpret “brilliant” as referring to “book,” the model must capture a long-range dependency - something RNNs often fail to maintain.

## The Attention Mechanism

Self-attention solves this by allowing each token to directly attend to all others in the sequence. Let  $X = [x_1, x_2, \dots, x_T]$ , where each  $x_i \in \mathbb{R}^d$  is the embedding of the  $i$ -th token.

Each input vector is linearly projected into three representations:

$$q_i = W_Q x_i \quad (11.1)$$

$$k_i = W_K x_i \quad (11.2)$$

$$v_i = W_V x_i \quad (11.3)$$

where:

- $q_i \in \mathbb{R}^{d_k}$  is the **Query vector** of token  $i$
- $k_i \in \mathbb{R}^{d_k}$  is the **Key vector** of token  $i$
- $v_i \in \mathbb{R}^{d_v}$  is the **Value vector** of token  $i$
- $W_Q, W_K, W_V \in \mathbb{R}^{d_k \times d}$  are learned projection matrices

The attention score between token  $i$  and token  $j$  is computed via scaled dot-product:

$$\alpha_{ij} = \frac{q_i^\top k_j}{\sqrt{d_k}} \quad (11.4)$$

Then, attention weights are obtained using the softmax function over the keys:

$$a_{ij} = \frac{\exp(\alpha_{ij})}{\sum_{j'=1}^T \exp(\alpha_{ij'})} \quad (11.5)$$

Finally, each token is updated by a weighted sum over the value vectors:

$$z_i = \sum_{j=1}^T a_{ij} v_j \quad (11.6)$$

This process creates a context-aware representation  $z_i$  for each token.

Let's consider for example the following self-attention map for the sentence "The solution that the engineer suggested was brilliant." Each row represents a word's query, and the columns represent the words it attends to. Darker cells indicate stronger attention weights. For example, the word "book" strongly attends to "brilliant", capturing the subject-verb agreement despite intervening tokens.

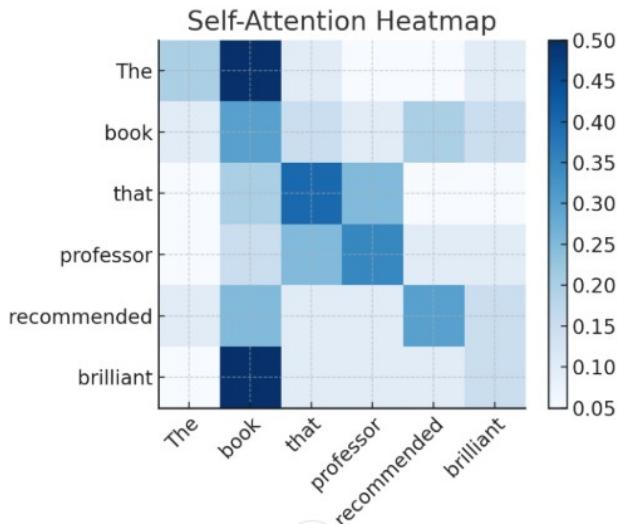


Figure 11.1: Self-Attention Heatmap for the sentence: "*The solution that the engineer suggested was brilliant.*"

## From Attention to the Transformer

A full Transformer model consists of an **encoder–decoder architecture**. The **encoder** maps an input sequence (e.g., a sentence) to a sequence of hidden representations using stacked layers of self-attention and feedforward components. The **decoder** uses masked self-attention (to prevent peeking ahead), cross-attention with the encoder output, and a feedforward network to generate the output tokens one by one.

This structure allows for conditional generation, making the Transformer a powerful **generative model** - suitable for tasks like translation, text generation, image captioning, and more.

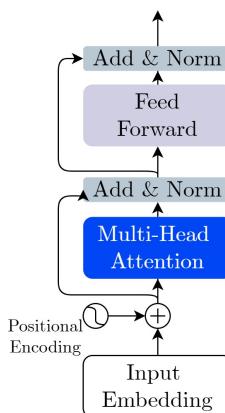


Figure 11.2: Encoder Illustration

The figure above illustrates the internal structure of a Transformer encoder block. At its core lies the **Multi-Head Attention** module - a parallelized extension of the attention mechanism - allowing the model to attend to different subspaces of the input representation simultaneously.

The attention output is then processed through a residual connection ( which we will explore more thoroughly in a later section) followed

by **Layer Normalization** (denoted as “Add & Norm”), which stabilizes the distribution of activations and facilitates deeper architectures.

Following the attention layer is a simple **position-wise feedforward network**, typically composed of two fully connected layers with a non-linearity (such as ReLU or GELU). Another residual connection and normalization step complete the block.

It is important to note that this architecture enables full permutation-aware contextual modeling without recurrence or convolution, relying entirely on learned attention patterns and positional encodings. This structure forms the backbone of modern language models and has been adapted to vision tasks (e.g., ViT), where image patches are treated as tokens in the same manner.

## Advantages and Generalization

Transformers offer several key advantages over traditional architectures: They allow full parallelization across sequence elements, which accelerates training. Self-attention mechanisms naturally model global context and long-range relationships, which RNNs struggle to capture. They scale well to very deep networks and support transfer learning through pretrained language models.

## Summary

Self-attention fundamentally changed the architecture of deep learning models. By allowing each token to selectively aggregate information from the entire sequence, it enabled powerful and scalable systems across domains. The Transformer architecture, with its encoder–decoder structure, serves as the foundation of today’s generative models in both language and vision.

## Lesson 12

### Normalization and Initialization

#### Introduction

In this lesson, we explore two critical methods in deep learning: **Normalization** and **Initialization**. These techniques are essential for stabilizing and accelerating the training of deep neural networks. Without proper normalization and initialization, models may converge slowly, get stuck in poor minima, or suffer from vanishing and exploding gradients. We begin by understanding why activations in deep networks must be normalized and then discuss how the initial choice of weights affects signal propagation.

#### Why Normalize?

During training, activations propagate through the network. However, as the depth increases, their distribution can change dramatically between layers or batches. This phenomenon, known as **internal covariate shift**, destabilizes learning and causes gradient magnitudes to fluctuate. To mitigate this, normalization techniques such as **Batch Normalization** and **Layer Normalization** are employed.

## Batch Normalization

Given a mini-batch of activations  $\{x_1, x_2, \dots, x_m\}$  from a particular neuron (or feature channel), Batch Normalization normalizes the values across the batch dimension:

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i \quad (12.1)$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad (12.2)$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}} \quad (12.3)$$

$$y_i = \gamma \hat{x}_i + \beta \quad (12.4)$$

Here:

- $\mu_B$  and  $\sigma_B^2$  are the mean and variance computed over the mini-batch,  $\hat{x}_i$  is the normalized activation,
- $\gamma$  and  $\beta$  are learnable parameters that allow rescaling and shifting,  $\varepsilon$  is a small constant added for numerical stability.

BatchNorm is typically inserted after the linear transformation and before the nonlinearity:

Linear → BatchNorm → ReLU

This improves gradient flow, allows for faster training, and enables the use of higher learning rates. However, it relies on sufficiently large batch sizes to produce stable statistics.

## Layer Normalization

In tasks involving very small batches or sequence-based architectures (such as NLP models or Transformers), BatchNorm may perform poorly due to unreliable batch statistics. In these cases, **Layer Normalization** is more effective, as it normalizes across the features of each individual sample - not across the batch.

Given an input vector  $x \in \mathbb{R}^d$  (e.g., the activations of one layer for a single sample), LayerNorm is computed as:

$$\mu = \frac{1}{d} \sum_{i=1}^d x_i \quad (12.5)$$

$$\sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2 \quad (12.6)$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \varepsilon}} \quad (12.7)$$

$$y_i = \gamma \hat{x}_i + \beta \quad (12.8)$$

Here, normalization is performed over all  $d$  features of the input vector for a single training example, and the same learnable scale and shift parameters ( $\gamma, \beta$ ) are applied.

This method is independent of batch size and works well in models where consistent behavior per-sample is needed, such as in recurrent networks or autoregressive decoders. This approach is especially effective in Transformers and recurrent networks where input patterns are highly variable and sequential.

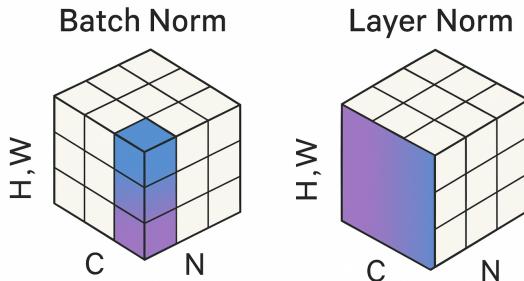


Figure 12.1: Comparison between Batch Norm (left) and Layer Norm (right). Batch Norm normalizes across the batch axis (N), while Layer Norm normalizes across all features within each sample (C, H, W).

## Weight Initialization

Before training begins, weights are typically initialized with random values. However, improper initialization can distort signal flow, cause gradient vanishing or explosion, and delay learning. The goal is to maintain controlled variance of both activations and gradients across all layers.

### Xavier Initialization

Designed for symmetric activations such as `tanh` or `sigmoid`, **Xavier Initialization** samples weights from a distribution with variance:

$$W_{ij} \sim \mathcal{N} \left( 0, \frac{2}{n_{\text{in}} + n_{\text{out}}} \right) \quad (12.9)$$

where  $n_{\text{in}}$  and  $n_{\text{out}}$  are the number of input and output units of the layer, respectively.

This balances the forward and backward signal flow to avoid amplification or decay.

## He Initialization

For non-symmetric activations like ReLU, which zero out negative inputs, **He Initialization** increases variance to:

$$W_{ij} \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}}}\right) \quad (12.10)$$

This ensures that the gradient flow is preserved despite sparsity introduced by ReLU.

## Summary

Deep networks are sensitive to the scale and distribution of activations and gradients. Normalization, particularly BatchNorm and LayerNorm, addresses instability by standardizing feature values, thereby enhancing the flow of information across layers. Meanwhile, proper initialization ensures the gradients are neither too large nor too small at the start of training.

Xavier initialization is ideal for networks using symmetric activations, while He initialization is tailored for ReLU-based models. These initialization schemes minimize signal distortion and improve learning efficiency.

Together, normalization and initialization form the backbone of stable, fast, and robust training dynamics in modern deep learning. In the next lesson, we will explore how to effectively augment training data, improving generalization without requiring the collection of additional samples.

## Lesson 13

### Data Augmentation

#### Introduction

A model may train stably, converge smoothly, and exhibit low training loss - yet fail dramatically in deployment. Why does this happen?

In the real world, data are noisy, inconsistent, and rarely resemble the perfectly curated samples we used in training. The key to generalization, the ability of the model to perform well on unseen inputs, is not always hidden in the network architecture. Often, it lies in the data itself. Data augmentation is the process of expanding and diversifying the dataset by applying label-preserving transformations. It increases robustness by simulating real-world variations, forcing the model to become invariant to irrelevant changes.

#### When Does Augmentation Help?

Data augmentation is especially helpful when the training dataset is small, imbalanced, or lacks sufficient diversity. In such cases, the model tends to overfit - achieving high accuracy on training data while failing to generalize to new or unseen examples. Augmentation introduces controlled variations into the data, helping the model become more

robust to changes in distribution. This is particularly important when the test domain includes patterns or distortions that are not well represented in the training set. A good augmentation strategy typically causes a slight increase in training loss, as the model faces harder and more varied inputs, but leads to a noticeable decrease in validation loss and better accuracy on out-of-distribution (OOD) data. As a result, augmentation serves as a simple yet powerful method to improve generalization and reliability, especially in vision tasks or when labeled data is limited.

## Augmentation vs. Generation

It is crucial to distinguish augmentation from generation. Augmentation does not fabricate new samples from scratch. Rather, it creates alternate views of the same sample while preserving its semantic label.

Let  $x$  be an input image and  $y$  its associated label. Augmentation constructs a transformed version  $x' = T(x)$  such that:

$$f(x') \approx f(x), \quad \text{and} \quad \text{label}(x') = y \quad (13.1)$$

Where  $T$  is a transformation function drawn from a distribution of admissible perturbations, and  $f$  is the model.

## Preserving Label Consistency

A fundamental requirement is that the transformed sample must retain its original label. Violating this can lead to noisy supervision and degraded performance.

For instance, rotating the digit "6" by  $180^\circ$  yields something visually indistinguishable from a "9" - yet labeling it as "6" would mislead the model. This is not augmentation; it is semantic corruption.

# Transformations

In visual domains, the following transformations are commonly applied:

- Horizontal Flip:  $x' = \text{Flip}_H(x)$  - suitable when object symmetry permits.
- Rotation:  $x' = \text{Rotate}_\theta(x), \theta \in [-15^\circ, 15^\circ]$
- Crop and Resize: Randomly cropping a sub-region and scaling back.
- Translation and Zoom: Small shifts and rescaling without altering class identity.
- Brightness and Contrast: Adjusting illumination levels.
- Color Jitter: Perturbing hue/saturation for color images.
- Additive Noise: Injecting Gaussian noise to simulate sensor variability.

## Practical Caution

Excessive augmentation can be counterproductive. If the resulting samples no longer resemble valid inputs - e.g., excessive distortion, loss of object structure - they may degrade learning. Visual inspection is essential.

## Implementation

A powerful open-source library for image augmentation is **Albumentations**, which integrates smoothly with both PyTorch and TensorFlow. This transformation pipeline randomly flips, rotates, brightens, and normalizes the image - all while preserving class semantics.

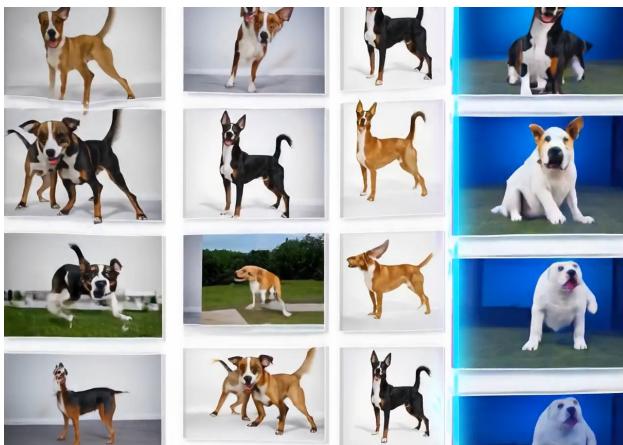


Figure 13.1: Augmentation Examples.

## Offline vs. Online Augmentation

There are two strategies for applying augmentation:

- **Offline augmentation:** Generates and stores new images in advance. It increases disk usage but offers fixed variability.
- **Online augmentation:** Applies random transformations in real-time during training. It introduces stochasticity and infinite diversity across epochs.

Most modern systems favor online augmentation due to its flexibility and memory efficiency.

## Best Practices

To maximize the effectiveness of augmentation:

1. **Validate visually:** Always inspect a batch of augmented images.
2. **Prioritize plausibility:** Avoid transformations that violate class identity.

3. **Prefer small, composable transformations:** Chain simple operations rather than aggressive single ones.
4. **Adapt per domain:** Augmentation for digits is not the same as for medical images or traffic signs.

## Summary

Data augmentation is one of the most powerful and accessible techniques to enhance model robustness. It simulates diversity, enforces invariance, and reduces reliance on handcrafted regularization. Unlike architecture changes or loss tweaks, augmentation transforms the data itself - helping the model encounter the world as it really is: noisy, variable, and unpredictable.

A well-designed augmentation pipeline does not just enrich the training set - it expands the model's perspective. It prepares the network not to memorize, but to understand.

## Lesson 14

### Advanced Optimization

Modern deep learning would not be possible without efficient and stable optimization algorithms. While neural networks may seem magical in their predictions, at their core they rely on careful and repeated updates to millions of parameters. In this lesson, we examine how these updates are performed - from the basic principle of gradient descent to the cutting-edge optimizers used these days: AdamW, Lion, and Adafactor. These optimizers determine how the model learns, how fast it converges, and whether it generalizes well.

### From Backpropagation to Optimization

Every training step begins with backpropagation, which calculates the gradients of the loss function with respect to each weight in the network. These gradients point in the direction that increases the loss the most. To improve the model, we want to go in the opposite direction (minimizing the loss). The optimizer is the algorithm that turns gradients into meaningful parameter updates.

## Gradient Descent: Back to Basic

As we saw previously, the classic update rule is simple and foundational:

$$\theta_{k+1} = \theta_k - \eta \cdot \nabla_{\theta} L(\theta_k) \quad (14.1)$$

Where  $\theta_k$  denotes the parameter vector at step  $k$ ,  $\eta$  is the learning rate, and  $\nabla_{\theta} L(\theta_k)$  is the gradient of the loss function  $L$  at that point. In intuitive terms, the model is sliding down the loss surface, guided by the local slope.

*Note: this is the exact same equation we encountered earlier, but in this context, the weight is denoted by  $\theta$  instead of  $w$ .*

Real-world training surfaces are not smooth hills; they are noisy, high-dimensional, and filled with ravines. Simple gradient descent struggles to deal with these complexities, which is why more advanced optimizers are required.

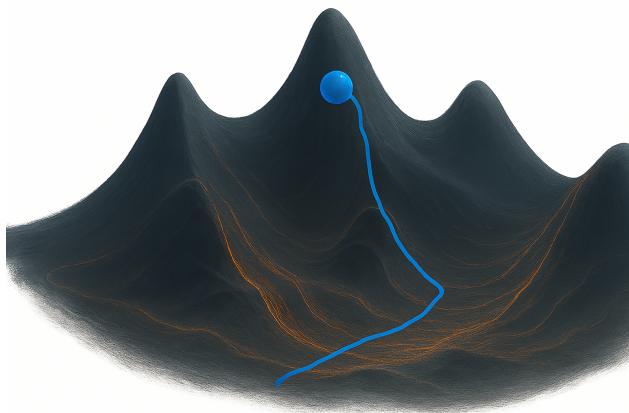


Figure 14.1: This image illustrates gradient descent: a ball rolls down a complex loss surface, guided by the steepest slope at each point. The path represents the optimization process as the model adjusts its parameters to minimize error. Peaks and valleys show the challenge of escaping local minima to reach the global minimum.

## Adam: Adaptive Moments

Adam, *Adaptive Moment Estimation*, improves upon plain gradient descent by combining two mechanisms: momentum and adaptive scaling. It keeps an exponentially decaying average of the gradients (the first moment) and their squares (the second moment):

$$m_k = \beta_1 \cdot m_{k-1} + (1 - \beta_1) \cdot g_k \quad (14.2)$$

$$v_k = \beta_2 \cdot v_{k-1} + (1 - \beta_2) \cdot g_k^2 \quad (14.3)$$

$$\theta_{k+1} = \theta_k - \eta \cdot \frac{m_k}{\sqrt{v_k} + \epsilon} \quad (14.4)$$

Where  $g_k$  is the gradient at step  $k$ , and  $m_k, v_k$  are moving averages of the gradient and squared gradient respectively. The term  $\epsilon$  (typically  $10^{-8}$ ) ensures numerical stability. Common default settings are  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ .

Adam adapts the step size individually for each parameter, making it robust and widely effective across tasks. It remains the default optimizer in many frameworks.

## AdamW: Decoupled Weight Decay

AdamW was introduced to correct a subtle flaw in Adam's implementation of regularization. In Adam, weight decay (L2 regularization) interacts undesirably with the adaptive step sizes. AdamW decouples the two, applying weight decay explicitly:

$$\theta_{k+1} = \theta_k - \eta \cdot \left( \frac{m_k}{\sqrt{v_k} + \epsilon} + \lambda \cdot \theta_k \right) \quad (14.5)$$

Here,  $\lambda$  is the weight decay coefficient (e.g., 0.01). This formulation

improves generalization and training stability, especially in large-scale transformer-based models. As a result, AdamW is the optimizer of choice for many state-of-the-art architectures.

## Lion: Lightweight Momentum with Signs

Lion, introduced in 2023 by researchers at Google, offers an elegant twist: instead of using the magnitude of the gradients or second moments, it updates weights based solely on the sign of the momentum:

$$m_k = \beta_1 \cdot m_{k-1} + (1 - \beta_1) \cdot g_k \quad (14.6)$$

$$\theta_{k+1} = \theta_k - \eta \cdot \text{sign}(m_k) \quad (14.7)$$

This simplification reduces memory usage and computational cost, while still preserving directionality. Lion is particularly effective for vision models like ViTs and has been shown to outperform AdamW in some training regimes, especially when GPU memory is constrained. Typical hyperparameters include:  $\eta = 3 \times 10^{-4}$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$ .

## Adafactor: Memory Efficiency at Scale

Adafactor was developed for extremely large models where even the memory required to store second-moment estimates becomes a bottleneck. Instead of maintaining full matrices for second-order moments, Adafactor factorizes them into row and column vectors, reducing the memory footprint from  $O(n^2)$  to  $O(n)$ . This makes it possible to train very large language models (like T5) on hardware with limited resources. Adafactor also supports an adaptive learning rate mechanism

based on parameter scale and can be used with scheduling techniques like cosine decay.

## Learning Rate Scheduling and Optimizers

Optimizers are only part of the picture; learning rate schedules play an equally critical role. A fixed learning rate rarely works well across all training stages. Common strategies include:

**Warmup + Decay:** Start with a low learning rate, increase gradually, then decay over time.

**Cosine Annealing:** Slowly reduce the learning rate following a cosine curve, often combined with warmup.

AdamW and Lion benefit significantly from cosine schedules, while Adafactor is often paired with adaptive schedules based on parameter norms. Selecting an appropriate learning rate schedule can greatly influence both convergence speed and final performance.

## Stability and Generalization Considerations

Adaptive optimizers like Adam and Lion often converge quickly, but may settle into sharp minima, points in parameter space that minimize training loss but generalize poorly. SGD with momentum, despite being slower, often finds flatter minima that lead to better test accuracy.

In very noisy settings, adaptive methods (especially AdamW) are often more stable. However, switching to simpler optimizers like SGD during fine-tuning can sometimes help improve robustness.

It is also important to note that different components of a model might benefit from different optimizers, for example, using AdamW for embeddings and SGD for classification heads.

## Feature Comparison Summary

Optimizer	Memory	Adapt. LR	Momentum	Use Case
SGD	Low	No	Yes	Robust tasks
AdamW	Medium	Yes	Yes	Most models
Lion	Low	No	Sign	ViTs, low RAM
Adafactor	Very Low	Yes	No	Huge LMs

Table 14.1: Modern optimizer summary

### Practical Tips

For noisy datasets or unstable gradients, AdamW is often a safe and effective choice. For resource-constrained GPU setups - especially when training ViTs - Lion may provide better throughput. When dealing with extremely large models (e.g.,  $>1\text{B}$  parameters), Adafactor becomes essential due to its efficiency.

### Summary

Modern optimizers go far beyond simple gradient descent. They integrate ideas from momentum, adaptive learning rates, memory efficiency, and regularization - each offering advantages in different scenarios.

AdamW is widely adopted for its balance of stability and performance; Lion offers simplicity and speed for vision models under resource constraints; Adafactor unlocks training at massive scale; and SGD, though basic, remains effective in low-noise or fine-tuning settings.

*An optimizer is not just a mathematical trick - it's your steering wheel. Choose it based on the terrain, the vehicle, and your destination.*

## Lesson 15

### Explainability: Understanding Model Decisions

#### Introduction

As deep learning models are increasingly used in high-stakes decision-making, the ability to explain their behavior becomes essential. In fields such as healthcare, finance, autonomous systems, and security, it is no longer sufficient for a model to be accurate. It must also be understandable.

Explainability, or interpretable machine learning, refers to the set of methods that provide insight into why a model produced a given output. These methods serve multiple purposes: they increase user trust, help developers debug unexpected behavior, ensure regulatory compliance, and uncover ethical or fairness concerns. Ultimately, a model that cannot be explained is a model that cannot be trusted.

#### Why Deep Learning Is Hard to Interpret

Linear models and decision trees offer transparency by construction: one can inspect coefficients or decision rules directly. Deep neural networks, in contrast, operate through multiple layers of non-linear transformations, often involving millions of parameters. The internal

representations they learn are distributed, entangled, and typically uninterpretable without post hoc analysis.

This complexity makes it difficult to determine what information influenced a particular decision, or to detect when the model is relying on spurious correlations or confounding variables.

## Case Study: Fraud Detection in Banking

Consider a bank deploying a deep learning model to detect fraudulent credit card transactions. During testing, the model achieves high accuracy, but the compliance team demands that each flagged transaction be accompanied by an explanation. When analysts investigate several false positives, they discover a troubling pattern: many are concentrated around a specific city.

Using SHAP explanations, they find that the “transaction location” feature has a high contribution score - but only when combined with certain merchant categories. This triggers a deeper audit, revealing that the model learned to associate “suspicious” activity with a combination of perfectly legitimate behaviors.

In a parallel effort, analysts apply Grad-CAM to a CNN-based submodule that analyzes scanned receipts. The heatmaps reveal that the model often focuses on highlighted text - not the transaction details - due to training data bias.

These two explanation methods, applied at different levels of the system, help the team retrain the model, improve the data pipeline, and satisfy regulatory demands for transparency.

## Grad-CAM: Visual Explanation in CNNs

In convolutional neural networks, a standard tool for visualizing model behavior is **Grad-CAM** (Gradient-weighted Class Activation Mapping). Grad-CAM generates a heatmap that highlights the spatial regions in the input image most responsible for the model’s decision.

Let  $y^c$  be the model’s score (pre-softmax) for class  $c$ , and let  $A^k \in \mathbb{R}^{H \times W}$  denote the  $k$ -th feature map from a selected convolutional layer of spatial dimension  $H \times W$ . The importance weight  $\alpha_k^c$  for feature map  $k$  is computed as:

$$\alpha_k^c = \frac{1}{Z} \sum_{i=1}^H \sum_{j=1}^W \frac{\partial y^c}{\partial A_{i,j}^k} \quad (15.1)$$

Here,  $Z = H \cdot W$  is the total number of spatial elements in  $A^k$ . These weights represent how strongly each feature map influences the class score.

The Grad-CAM heatmap is then computed as:

$$L_c^{\text{Grad-CAM}} = \text{ReLU} \left( \sum_k \alpha_k^c A^k \right) \quad (15.2)$$

The ReLU operation ensures that only the positively contributing regions are retained, focusing on evidence in favor of class  $c$ . This heatmap is upsampled to the input resolution for visualization.

## Global Importance via Permutation-Based Analysis

In tabular models, interpretability often focuses on identifying which features are globally important. While tree-based models such as XG-Boost provide internal metrics (e.g., gain, coverage), a model-agnostic

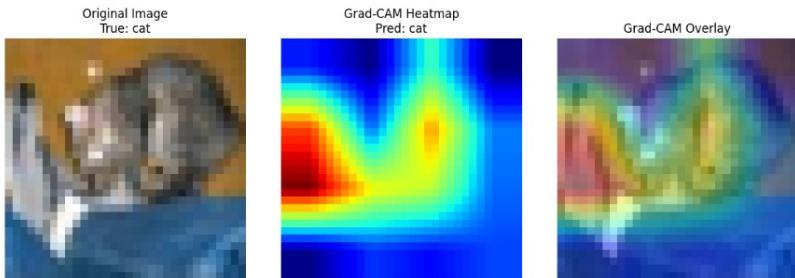


Figure 15.1: Example - visualizing a cat prediction in CIFAR-10. The image shows a Grad-CAM explanation for a CNN trained on CIFAR-10. The model correctly classifies the image as a *cat*, and Grad-CAM highlights the most relevant regions that led to this prediction.

and intuitive approach is **permutation importance**.

Let  $f$  be a trained model and  $\mathcal{D}$  a validation dataset. For each feature  $x_i$ , the permutation importance is computed by shuffling its values across examples and measuring the change in performance. Formally, let  $M_0$  be the baseline metric (e.g., accuracy) and  $M_i$  be the metric after permuting  $x_i$ :

$$\Delta_i = M_0 - M_i \quad (15.3)$$

A large  $\Delta_i$  indicates that feature  $x_i$  is critical to performance, while a small value suggests limited influence.

## SHAP: Game-Theoretic Local Explanations

SHAP (SHapley Additive exPlanations) is a powerful method for attributing a model's output to its input features. It is grounded in Shapley values from cooperative game theory and satisfies four key axioms: efficiency, symmetry, dummy, and additivity - which ensure fairness and consistency in feature attribution.

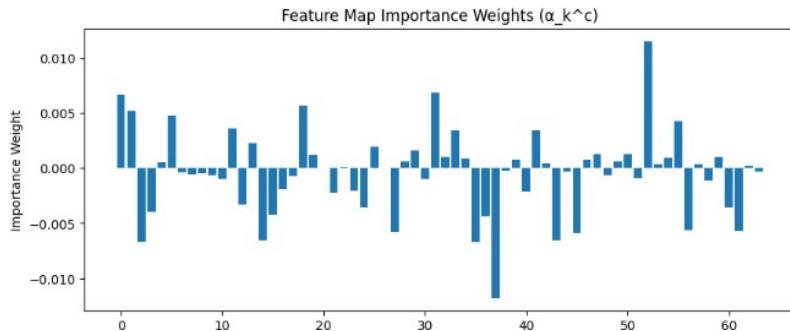


Figure 15.2: Bar chart showing the Grad-CAM importance weights  $\alpha_k^c$  for each feature map. Each feature map captures a specific visual pattern learned by the CNN (e.g., textures, shapes, parts). Positive weights indicate maps that supported the "cat" prediction; negatives were filtered by ReLU.

Let  $F$  denote the full set of input features, and let  $f(S)$  be the model's prediction using only the feature subset  $S \subseteq F$ . The SHAP value for feature  $i$  is defined as:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f(S \cup \{i\}) - f(S)] \quad (15.4)$$

In this equation:

- $\phi_i$  is the attribution score for feature  $i$ .
- $f(S)$  is the expected model output using subset  $S$  of features.
- The combinatorial factor weights each subset proportionally to its likelihood in random feature orderings.

This formulation considers all possible coalitions of features and measures the marginal contribution of  $x_i$  to the model's output.

In practice, computing Eq. 15.4 exactly is computationally prohibitive for large  $F$ . Approximations such as Kernel SHAP (model-agnostic) or Deep SHAP (for neural networks) are used instead. These methods approximate  $f(S)$  using background distributions or reference samples.

Despite its cost, SHAP provides instance-level explanations that are robust, additive, and intuitive, making it one of the most widely used tools in regulated industries.

## Text Explainability via Attention Weights

In transformer-based models such as BERT or GPT, interpretability is embedded into the architecture via the self-attention mechanism. Each input token computes a weighted average over all other tokens, producing a contextualized representation.

Let's recall the attention formula from earlier lessons. Given a sequence of input embeddings, the attention weight between token  $i$  and token  $j$  is defined as:

$$\text{Attention}(i, j) = \frac{\exp(q_i^\top k_j / \sqrt{d_k})}{\sum_{j'=1}^T \exp(q_i^\top k_{j'} / \sqrt{d_k})} \quad (15.5)$$

Here:

- $q_i, k_j \in \mathbb{R}^{d_k}$  are the query and key vectors for tokens  $i$  and  $j$ , respectively.
- $d_k$  is the dimensionality of the key/query space (used for scaling).
- $T$  is the total number of tokens in the sequence.

These attention weights form a matrix  $A \in \mathbb{R}^{T \times T}$  that can be visualized to reveal which words influence each other in a given layer and

head. While attention does not directly explain predictions, it provides insight into the model’s internal reasoning, particularly in co-reference, translation, or syntax-sensitive tasks.

## Limitations and Interpretability Tradeoffs

Despite their utility, explainability methods are not without caveats. Grad-CAM is limited by spatial resolution and may produce noisy heatmaps, especially in deeper layers. SHAP, though principled, can be unstable under feature correlation and is expensive for high-dimensional inputs. Attention weights reflect association, not causality, and may vary across heads and layers without clear interpretive alignment.

Explanations can also be misleading. A model might assign high SHAP values to features that are proxies for deeper confounders, or attention heads may highlight syntactic tokens with no semantic relevance.

Thus, explainability should be used as a diagnostic lens - not a ground truth oracle. Combining these methods with domain knowledge, robustness checks, and human evaluation is essential.

## Summary

Explainability transforms deep learning from a black box to an auditable, actionable system. Grad-CAM highlights influential image regions in CNNs. SHAP offers mathematically rigorous attributions for individual predictions, rooted in game theory. Permutation-based techniques offer simple but effective global insights, and attention weights open windows into token-level interactions in Transformers.

Each technique has strengths and tradeoffs. Their selection depends on the data modality, model type, and interpretability goals.

## Lesson 16

### TensorFlow vs. PyTorch

#### Introduction

After learning how to design, train, and interpret neural networks, one practical question often arises: *which framework should we use?* The two most popular libraries today are TensorFlow and PyTorch. Both are widely used, well-maintained, and capable of supporting both research and production environments. They share much of the same functionality - yet differ in philosophy, syntax, and typical usage patterns.

In this lesson, we compare the two frameworks through structure, style, and ecosystem. We also see a hands-on classification example using Fashion MNIST to highlight the real differences in code and abstraction level.

#### Historical Background and Ecosystem

TensorFlow was introduced by Google in 2015 as a scalable, production-oriented platform. It became the default framework for many enterprise applications, thanks to tools like TensorFlow Serving, TensorFlow Lite, and TensorFlow.js.

PyTorch was released by Facebook in 2016 and rapidly gained popularity in the research community. Its native Python syntax, dynamic execution model, and flexibility made it the preferred choice for experimentation and rapid prototyping.

Over time, both frameworks have evolved. TensorFlow 2 adopted eager execution by default, and PyTorch now includes tools for deployment, such as TorchServe and TorchScript. Both support ONNX (Open Neural Network Exchange) and integrate well with modern ML tools.

## Execution Models: Static vs. Dynamic

Historically, TensorFlow required users to first define a static computation graph, which was then executed inside a session. This model made optimization easier but reduced flexibility.

PyTorch, from the start, used dynamic computation: code is executed line-by-line as written in Python. This makes it easier to debug and to build models with complex control flow.

In TensorFlow 2, eager execution is now default. However, it also provides the `@tf.function` decorator to convert Python functions into graph representations - combining flexibility with performance when needed.

## Case Study: Fashion MNIST Classification

We now demonstrate both frameworks on the same classification task: Fashion MNIST. The dataset includes grayscale images ( $28 \times 28$  pixels) of clothing items, divided into 10 classes. The goal is to implement a simple neural network for classification.

## PyTorch Implementation

Listing 16.1: PyTorch implementation of a simple classifier

```
import torch
import torch.nn as nn
from torchvision import datasets, transforms
from torch.utils.data import DataLoader

train_data = datasets.FashionMNIST('.', train=True,
    download=True, transform=transforms.ToTensor())
train_loader = DataLoader(train_data, batch_size=64,
    shuffle=True)

class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.net = nn.Sequential(
            nn.Flatten(),
            nn.Linear(784, 128),
            nn.ReLU(),
            nn.Linear(128, 10)
        )

    def forward(self, x):
        return self.net(x)

model = Net()
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
loss_fn = nn.CrossEntropyLoss()

for x, y in train_loader:
    pred = model(x)
    loss = loss_fn(pred, y)
    loss.backward()
    optimizer.step()
    optimizer.zero_grad()
```

This structure requires the user to explicitly define the forward pass and implement each training step. It offers full control over the learning process.

## TensorFlow Implementation

Listing 16.2: TensorFlow (Keras) implementation of the same model

```
import tensorflow as tf

(train_x, train_y), _ = tf.keras.datasets.fashion_mnist.
    load_data()
train_x = train_x / 255.0

model = tf.keras.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(10)
])

model.compile(optimizer='adam',
    loss=tf.keras.losses.SparseCategoricalCrossentropy(
        from_logits=True),
    metrics=['accuracy'])

model.fit(train_x, train_y, batch_size=64, epochs=5)
```

TensorFlow, through Keras, offers a more declarative interface. It abstracts the training loop and internal logic, making it ideal for quick experimentation or onboarding new developers.

## Syntax, Modularity, and Abstraction

The code above illustrates the core stylistic difference:

- PyTorch emphasizes transparency and modularity. The user defines the model, forward pass, and optimizer logic explicitly.

- TensorFlow with Keras abstracts these details, enabling fast development but requiring deeper configuration for advanced tasks.

PyTorch code generally feels like native Python, while TensorFlow code - particularly when defining functional models or using graph mode - may feel more like a framework with its own conventions.

## Ecosystem and Production Readiness

When selecting a framework, consider not only development but also deployment:

TensorFlow supports:

- **TensorFlow Serving** for REST and gRPC APIs
- **TensorFlow Lite** for mobile/embedded inference
- **TensorFlow.js** for browser-based execution

PyTorch supports:

- **TorchScript** for exporting models to serialized graphs
- **TorchServe** for serving models
- Export to **ONNX**, enabling deployment in other frameworks

Both libraries support mixed-precision training, GPU/TPU acceleration, and integration with Hugging Face, Weights & Biases, and MLflow.

## When to Use Each Framework

The choice depends on context. PyTorch is excellent for rapid experimentation, research, and teaching. It provides fine-grained control and clear debugging. TensorFlow excels in end-to-end pipelines, long-term production systems, and mobile/web deployments.

Because of ONNX support, models can be ported between them. It is increasingly common to prototype in PyTorch and convert to TensorFlow for production - or vice versa.

## Summary

Both TensorFlow and PyTorch are capable, modern deep learning frameworks. They share many features and are converging in design. The key distinctions lie in their syntax, abstraction level, and ecosystem tools. Use PyTorch for control, flexibility, and dynamic modeling. Use TensorFlow (Keras) for rapid development and streamlined production deployment.

## Lesson 17

### Working Effectively with Google Colab

#### Introduction

Google Colab has become the starting point for many developers, students, and researchers in deep learning - and for good reason. It is free, cloud-hosted, supports GPUs, and allows seamless sharing and collaboration. But like any cloud-based environment, effective use requires understanding its limitations, shortcuts, and hidden capabilities.

This lesson provides a workflow for working efficiently in Google Colab, including hardware configuration, file management, installation tips, model training, and best practices to avoid timeouts and maximize productivity.

#### Getting Started with the Colab Environment

When you open a new Colab notebook, you're working on a cloud-hosted Jupyter Notebook running on Google's infrastructure. Each notebook cell can include code, Markdown explanations, or results. Code is executed in an isolated virtual machine, and resources such as memory, disk, and GPU usage are shown at the top right of the interface.

You can toggle between code and text blocks using the "+ Code" and "+ Text" buttons on the top bar.

```

import tensorflow as tf
import numpy as np
import matplotlib.pyplot as plt
from tensorflow.keras import layers, models
import cv2

# Load and prepare data
print("Loading CIFAR-10 dataset...")
(x_train, y_train), (x_test, y_test) = tf.keras.datasets.cifar10.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

class_names = ['airplane', 'automobile', 'bird', 'cat', 'deer',
               'dog', 'frog', 'horse', 'ship', 'truck']

# Create and train model using Functional API
print("Creating simple CNN...")
inputs = tf.keras.Input(shape=(32, 32, 3))
x = layers.Conv2D(32, (3, 3), activation='relu')(inputs)
x = layers.MaxPooling2D((2, 2))(x)
x = layers.Conv2D(64, (3, 3), activation='relu')(x)
x = layers.MaxPooling2D((2, 2))(x)
conv_output = layers.Conv2D(64, (3, 3), activation='relu', name='target_conv')(x)

```

Figure 17.1: Google Colab Screenshot.

## Enabling GPU or TPU Acceleration

Colab allows you to request access to a GPU or TPU, which significantly accelerates training.

1. Go to Runtime > Change runtime type
2. Select GPU or TPU
3. Click Save

Once activated, you can confirm GPU access via:

```

import torch
print(torch.cuda.is_available())

```

or by inspecting the current GPU device:

```
!nvidia-smi
```

A common output includes devices like Tesla T4 or V100.

## Installing Python Packages

Colab supports direct installation of Python packages using pip inside notebook cells. For example:

```
!pip install transformers --upgrade --quiet
```

The `--quiet` flag hides lengthy installation logs.

## Connecting to Google Drive

To save models, outputs, or datasets across sessions, mount your Google Drive:

```
from google.colab import drive  
drive.mount('/content/drive')
```

Files are then accessible via `/content/drive/MyDrive/`. This ensures that important results are not lost when the runtime resets.

## Uploading Files Manually

You can upload files from your local machine either by:

- Dragging them into the "Files" pane
- Using the following command:

```
from google.colab import files  
files.upload()
```

Uploaded files will be stored temporarily in the notebook's local filesystem.

## Sharing and Collaborating

Each Colab notebook can be shared like any Google Doc. Click the “Share” button at the top-right, and choose access permissions. You can also open Colab notebooks directly from GitHub.

## Training a Neural Network with GPU

Let’s train a simple network on Fashion MNIST using PyTorch with GPU support:

```
import torch
import torchvision
import torchvision.transforms as transforms
from torch import nn, optim
from tqdm import tqdm

device = 'cuda' if torch.cuda.is_available() else 'cpu'

train_data = torchvision.datasets.FashionMNIST(root='.',
    train=True, download=True,
        transform=transforms.ToTensor())
train_loader = torch.utils.data.DataLoader(train_data,
    batch_size=64, shuffle=True)

model = nn.Sequential(
    nn.Flatten(),
    nn.Linear(784, 128),
    nn.ReLU(),
    nn.Linear(128, 10)
).to(device)

loss_fn = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters())

for epoch in range(3):
```

```
total_loss = 0
for x, y in tqdm(train_loader):
    x, y = x.to(device), y.to(device)
    pred = model(x)
    loss = loss_fn(pred, y)
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
    total_loss += loss.item()
print(f"Epoch {epoch+1}, Loss: {total_loss:.3f}")
```

## Monitoring with TensorBoard

To monitor training visually, Colab supports TensorBoard. First, load the extension:

```
%load_ext tensorboard
%tensorboard --logdir runs
```

Inside your training loop, log metrics using:

```
from torch.utils.tensorboard import SummaryWriter
writer = SummaryWriter()

# During training
writer.add_scalar('Loss/train', loss.item(), step)
```

## Advanced Tips for Smooth Usage

- Save model checkpoints:

```
torch.save(model.state_dict(), 'model.pth')
```

- Check disk usage:

```
!du -sh /content
```

- List files:

```
!ls -lh
```

- Measure execution time of a cell:

```
%%time
```

- View memory usage:

```
!free -h
```

## Avoiding Disconnections and Timeouts

Google Colab automatically disconnects sessions after 90 minutes of inactivity. To avoid losing progress, it's recommended to periodically run a cell to keep the session active, save models and logs directly to Google Drive, and divide long training processes into shorter runs. For extended runtime and improved hardware, consider upgrading to Colab Pro or Pro+.

## Colab Gemini Assistant (Pro+ Only)

In Colab Pro+, a Gemini-based assistant appears during programming. It can explain code you've written, suggest completions or improvements, and answer natural language queries about ML workflows.

## Summary

Google Colab offers a versatile, accessible, and GPU-accelerated platform for deep learning. With the right workflow, it becomes much more than a scratchpad - it becomes a full lab in the cloud. Learn its shortcuts, use Drive wisely, monitor training visually, and work modularly.

## Lesson 18

### Mixed Precision Training

#### Introduction

By now, we've seen how deep networks are built, trained, optimized, and evaluated. But there's one final layer that affects all of this - a hidden dimension where time, memory, and precision intersect.

Mixed Precision Training is not a new optimizer, a regularization trick, or a new loss function. It's a fundamental change in how we represent and compute the numbers inside our networks - and, if used correctly, it can lead to faster training, lower memory consumption, and sometimes even better generalization.

#### What Precision Means in Practice

Most neural networks today are trained using 32-bit floating point numbers (known as FP32). This is the default format for storing weights, computing activations, and tracking gradients. It offers excellent numerical stability, but comes with a cost: every number requires 32 bits, and every operation uses a relatively expensive data path.

Modern GPUs, particularly those from NVIDIA's Turing and Ampere families, support an alternative: 16-bit floating point numbers

(FP16). These are smaller, faster, and more memory-efficient - but they sacrifice a significant amount of numerical precision. The difference is not just in storage. FP32 uses 8 bits for the exponent and 23 for the mantissa. FP16, in contrast, uses just 5 exponent bits and 10 for the mantissa. This narrower range can lead to instability if not handled carefully.

Despite this, using FP16 can double memory throughput and reduce training time significantly. The challenge is deciding where it's safe to use it - and where it's not.

## The Mixed Precision Approach

Instead of choosing between FP32 and FP16, modern training pipelines combine both. This is the idea behind mixed precision. By default, calculations that are robust to low precision are carried out in FP16. Those that are sensitive - especially those involving loss values, small gradients, or weight updates - are done in full precision.

This separation allows us to leverage the speed and memory benefits of FP16 without compromising the model's stability. The transition is handled automatically by frameworks like PyTorch and TensorFlow, using tools such as AMP - Automatic Mixed Precision.

## AMP in PyTorch: A Minimal Example

PyTorch makes mixed precision remarkably easy to implement. The key components are two utilities from `torch.cuda.amp`: `autocast`, which handles the automatic casting of operations to FP16 where appropriate, and `GradScaler`, which protects against gradient underflow by scaling the loss before backpropagation and unscaling it before the optimizer step.

Here is a full example of training DNN on CIFAR-10 using AMP:

Listing 18.1: Training with AMP: ResNet18 on CIFAR-10

```
import torch
from torch import nn, optim
from torchvision.models import resnet18
from torchvision.datasets import CIFAR10
from torchvision.transforms import ToTensor
from torch.utils.data import DataLoader
from torch.cuda.amp import autocast, GradScaler
from tqdm import tqdm

device = 'cuda' if torch.cuda.is_available() else 'cpu'

train_data = CIFAR10(root='.', train=True, download=True,
    transform=ToTensor())
train_loader = DataLoader(train_data, batch_size=64,
    shuffle=True)

model = resnet18(num_classes=10).to(device)
optimizer = optim.Adam(model.parameters())
loss_fn = nn.CrossEntropyLoss()
scaler = GradScaler()

for epoch in range(3):
    total_loss = 0
    for x, y in tqdm(train_loader):
        x, y = x.to(device), y.to(device)
        optimizer.zero_grad()
        with autocast():
            pred = model(x)
            loss = loss_fn(pred, y)
            scaler.scale(loss).backward()
            scaler.step(optimizer)
            scaler.update()
            total_loss += loss.item()
    print(f"Epoch {epoch+1}, Loss: {total_loss:.3f}")
```

This version of the training loop requires only three changes: wrapping the forward pass in `autocast`, replacing `loss.backward()` with a scaled backward pass, and using `scaler.step` instead of the standard optimizer step. These small adjustments can yield dramatic speedups - especially when training large models or using large batch sizes.

## How to Know If It's Helping

To measure the real benefit of mixed precision, start by tracking runtime. You can use `%%time` to profile a training cell or compare epochs side-by-side. For memory, run `!nvidia-smi` in a separate cell and observe the reduction in consumption. One of the most immediate benefits is the ability to increase batch size without running into memory limits - often doubling or tripling your throughput.

## When It Works, and When to Watch Out

Mixed precision training works astonishingly well in most scenarios - but not all. Models with very deep or wide layers, such as Transformers, benefit the most. So do tasks with large inputs, like image generation or sequence modeling. In these cases, AMP should be considered default.

However, there are exceptions. If you're training on a GPU without Tensor Cores - such as older Kepler or Maxwell cards - you may see no benefit, or even a slowdown. Likewise, if your task involves detecting subtle statistical anomalies, such as medical diagnostics or rare-event classification, the reduced precision **may harm** the model's sensitivity. Finally, some edge cases produce NaNs during training due to underflow or exploding gradients. When this happens, check that you're using `GradScaler`. If issues persist, lower the learning rate or temporarily revert to FP32.

## A New Default in 2025

The deep learning ecosystem in 2025 assumes mixed precision. Frameworks like PyTorch Lightning and TensorFlow default to AMP. Most Kaggle notebooks use it out-of-the-box. Many Hugging Face training scripts enable it by default. Unless you're in a research setting that explicitly requires high numerical fidelity, there is no longer a reason not to use AMP.

If you're training a large model, using a modern GPU, or working with batch sizes above 128 - you should be using mixed precision. The benefits are tangible: more throughput, faster convergence, and fewer memory bottlenecks.

## Summary

Mixed Precision Training is one of the rare techniques in deep learning that improves both performance and efficiency without requiring architectural changes. It allows models to run faster and consume less memory while still achieving the same - or better - results. Thanks to tools like AMP, integrating it into your pipeline is trivial.

The core idea is simple: compute quickly where you can, compute carefully where you must. Like all good engineering, it's about balance. And in the era of large models and limited budgets, that balance can make all the difference.

*It's not enough to train accurately - you must also train efficiently.*

## Lesson 19

### Transfer Learning and Fine-Tuning

#### Introduction

Training a deep learning model from scratch often requires massive datasets, significant compute resources, and extensive tuning. However, in practice, it is often more effective to begin not from scratch, but from a pre-trained model—a model that has already learned useful representations from large-scale data. This approach is called **transfer learning**.

Transfer learning leverages the feature extraction capabilities of a model trained on a general domain (e.g., ImageNet or Wikipedia) and adapts it to a specific downstream task (e.g., classifying medical images or analyzing sentiment). The process of refining this pre-trained model to fit the new task is known as **fine-tuning**.

In this lesson, we will explore the motivation, strategies, and implementation of transfer learning and fine-tuning, with examples from computer vision and NLP.

## What is Transfer Learning?

In classical supervised learning, we train a model from randomly initialized weights using labeled examples. In transfer learning, we *start* from a model whose weights were trained on a related-but much larger- dataset.

In deep networks, early layers tend to learn general-purpose features (edges, textures, local patterns), while deeper layers capture more domain-specific representations. Transfer learning assumes that these lower-level representations are reusable across tasks.

*Example:* A convolutional neural network (CNN) trained to classify objects in natural images (e.g., dogs, trucks) can be adapted to classify product photos or X-ray scans by reusing the early layers and updating only the final classification layers.

## Common Pretrained Models

- **Computer Vision (CV):** ResNet, MobileNet, EfficientNet, Vision Transformers (ViT) – typically trained on ImageNet.
- **Natural Language Processing (NLP):** BERT, RoBERTa, GPT, DistilBERT – typically trained on large textual corpora.
- **Audio / Speech:** wav2vec, Whisper – trained on raw waveforms or spectrograms.

Pretrained models are available in public repositories such as Hugging Face Model Hub or Torchvision, and can be loaded with a single line of code.

## Three Fine-Tuning Strategies

Transfer learning is typically implemented via one of the following:

1. **Feature Extraction (Fixed Base)**: Freeze all pre-trained layers and train only a new output layer on the new task.
2. **Partial Fine-Tuning**: Freeze early layers, fine-tune later layers and the output head.
3. **Full Fine-Tuning**: Allow all weights to be updated, usually with a lower learning rate.

Each approach balances flexibility and stability. Full fine-tuning can yield higher performance on domain-specific tasks, but requires careful regularization.

## Practical Example: Fine-Tuning ResNet on Fashion MNIST

Suppose we want to classify grayscale clothing images in Fashion MNIST using transfer learning.

### Step 1: Load Pretrained Model

```
from torchvision.models import resnet18  
model = resnet18(pretrained=True)
```

### Step 2: Freeze All Layers

```
for param in model.parameters():  
    param.requires_grad = False
```

## Step 3: Replace Output Layer

```
import torch.nn as nn
model.fc = nn.Linear(512, 10) # 10 classes in Fashion
MNIST
```

## Step 4: Train the Output Head

Use a small learning rate (e.g.,  $10^{-4}$ ) and train the new classification layer. Optionally, unfreeze additional layers progressively.

## Mathematical Considerations

During fine-tuning, it is essential to use an appropriate learning rate  $\eta$ . Let  $\theta$  denote the parameter vector of the pre-trained model and  $\mathcal{L}$  the task-specific loss. The update step is:

$$\theta_{k+1} = \theta_k - \eta \cdot \nabla_{\theta} L(\theta_k) \quad (19.1)$$

When fine-tuning, we typically apply a smaller  $\eta$  for the pre-trained layers and a larger  $\eta$  for newly initialized layers. This prevents catastrophic forgetting of the previously learned representations.

## Fine-Tuning in NLP

Pretrained language models like BERT are adapted using a similar approach. For example:

```
from transformers import BertForSequenceClassification

model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased", num_labels=2)
```

This loads BERT with a classification head suited for binary sentiment analysis. Fine-tuning proceeds via standard backpropagation on labeled text.

## When Not to Use Transfer Learning

While powerful, transfer learning has limitations:

- If the source and target domains are highly dissimilar (e.g., dogs → medical scans)
- If the pretrained model introduces bias irrelevant to the new task
- If you have sufficient domain-specific data and compute to train from scratch

Always compare against a baseline trained from scratch.

## Best Practices

- Always start with `requires_grad = False` and unfreeze gradually.
- Use a lower learning rate for pretrained layers.
- Monitor performance on validation to detect overfitting.
- Consider layer-wise learning rate scheduling for large models.

## Summary

Transfer learning allows us to build accurate models quickly, with less data and fewer resources. It works by reusing general representations from a base model trained on a related task, and adapting them via

fine-tuning. This approach is now standard in nearly every subdomain of deep learning.

Whether classifying images or analyzing text, transfer learning transforms model development from training to adaptation, accelerating development and improving generalization.

## Lesson 20

### Saving, Loading, and Versioning Models

#### Introduction

After investing hours - or days - in training a deep learning model, the last thing you want is to lose the result. Accuracy, after all, means nothing if the model cannot be reused, reproduced, or deployed.

This lesson covers the professional practice of saving, loading, and versioning models in a way that supports collaboration, reproducibility, and long-term maintainability. We will explore the differences between formats like TorchScript, SavedModel, and ONNX, and walk through a working example using a CNN trained on Fashion MNIST.

#### Why Saving Models Is More Than Technical

Saving a model is not just a technical necessity - it's a matter of engineering discipline. If you don't preserve the trained parameters, the model structure, and the environment in which it was created, you've essentially lost the experiment. Even a 94% accuracy result is meaningless if you can't reproduce it.

In real-world scenarios, failure to save models and log configurations can lead to data loss, wasted time, and broken downstream pipelines.

Good model saving is a commitment to your future self - or to the next person who picks up your code.

## The Three Primary Saving Formats

Across the deep learning ecosystem, there are three dominant formats used to export and preserve trained models.

### TorchScript (PyTorch)

TorchScript is PyTorch's solution for exporting models into a standalone, serializable format. It compiles a model into a static computation graph, allowing it to run outside the Python runtime - for example, in C++ applications or on mobile devices.

Listing 20.1: Saving a TorchScript model

```
traced = torch.jit.trace(model, example_input)
traced.save("model.pt")
```

To load and use the model for inference:

Listing 20.2: Loading and evaluating TorchScript

```
model = torch.jit.load("model.pt")
model.eval()
```

This format is compact, fast, and highly portable - ideal for integrating into production services with minimal dependencies.

### SavedModel (TensorFlow)

In TensorFlow, the default export format is SavedModel. It preserves not only the weights and computation graph, but also function signatures and metadata. Saving is straightforward:

Listing 20.3: Saving a TensorFlow model

```
model.save("my_model")
```

Loading is just as simple:

Listing 20.4: Loading a SavedModel

```
loaded = tf.keras.models.load_model("my_model")
```

This format is particularly suitable for deployment via TensorFlow Serving, TensorFlow Lite, or Google Cloud AI Platform.

## ONNX (Open Neural Network Exchange)

ONNX is an open, cross-framework format designed to support interoperability. A model trained in PyTorch can be exported to ONNX and then used in other frameworks, runtimes, or languages.

Listing 20.5: Exporting to ONNX from PyTorch

```
torch.onnx.export(model, dummy_input, "model.onnx",
                  input_names=["input"], output_names=["output"])
```

ONNX is ideal for serving models in environments that are not tied to a specific deep learning library, such as C++, Java, or REST APIs powered by ONNX Runtime or Triton Inference Server.

## Practical Example: CNN on Fashion MNIST

Let's walk through a concrete example. We define and train a small convolutional network for image classification, and then save it in three formats: raw weights, TorchScript, and ONNX.

Listing 20.6: CNN definition in PyTorch

```

import torch
import torch.nn as nn

class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.net = nn.Sequential(
            nn.Conv2d(1, 16, 3, padding=1),
            nn.ReLU(),
            nn.Flatten(),
            nn.Linear(16 * 28 * 28, 10)
        )
    def forward(self, x):
        return self.net(x)

model = Net()
# Assume model is already trained

```

First, we can save just the weights:

Listing 20.7: Saving and loading state dict

```

torch.save(model.state_dict(), "cnn_weights.pth")

# To reload:
model.load_state_dict(torch.load("cnn_weights.pth"))
model.eval()

```

Then, for deployment outside PyTorch:

Listing 20.8: Saving as TorchScript

```

example_input = torch.randn(1, 1, 28, 28)
traced = torch.jit.trace(model, example_input)
traced.save("cnn_scripted.pt")

```

And finally, to export as ONNX:

Listing 20.9: Exporting CNN to ONNX

```
dummy_input = torch.randn(1, 1, 28, 28)
torch.onnx.export(model, dummy_input, "cnn.onnx",
                  input_names=["input"], output_names=["
output"])
```

## Choosing the Right Format

The best format depends on your deployment context, ecosystem, and goals. Each option has specific advantages, as summarized below:

Format	Primary Use	Advantages
TorchScript	Fast inference, C++/mobile integration	Lightweight, portable, runs without Python
SavedModel	TensorFlow Serving or TFLite	Includes graph, weights, and signatures; cloud-friendly
ONNX	Interoperability between frameworks	Open format, supported by multiple runtimes (Triton, ONNX Runtime, Optimum)

Table 20.1: Comparison of model export formats

## Versioning Like a Professional

Model saving is only half the story. Serious projects require versioning. Every saved model should be stored with a clear versioned filename, such as `model_v2.1.onnx`. Also, it may contain a companion log or metadata file that includes: training date, a short description of changes, key hyperparameters, final evaluation metrics on the test set. This information ensures that you - or anyone else - can reproduce the results months later. Consider integrating with tools like Git,

MLflow, or DVC to track experiments and model artifacts across time and machines.

## Summary

Saving a model is not the end of a training run - it's the start of its real-world life. Whether you're deploying to production, sharing with a team, or revisiting an old experiment, proper model serialization and versioning are essential.

Formats like TorchScript, SavedModel, and ONNX make it easy to move models across languages, runtimes, and infrastructures. But no format can preserve intent or context - that part is up to you.

*Don't just build a great model - make sure it's saved, understood, and ready for the future.*

## Lesson 21

### Basic Industrial Deployment

#### Introduction

After training and saving a deep learning model, the next essential step is to make it usable - not just from within your notebook, but by other systems. This is the bridge between modeling and production. It's how an image classifier becomes a real-time service, how a recommendation engine powers a web app, and how deep learning becomes part of a working product.

In this lesson, we will show how to deploy a trained image classification model as a RESTful API using FastAPI and TorchScript. The goal: receive an image via HTTP, return a prediction in JSON - clean, fast, and production-ready.

#### From Model to Service

In a production setting, we don't re-run training scripts or notebooks to get predictions. Instead, we wrap the model in a persistent service - one that listens for incoming requests, processes them, and responds with predictions. This architecture allows the model to be integrated into websites, mobile apps, dashboards, IoT systems, or any other software

interface.

A model API receives input (like an image), preprocesses it, performs inference using a preloaded model, and returns the result in a structured format such as JSON.

## FastAPI: Serving Inference with Clarity

FastAPI is a Python web framework designed for building high-performance APIs. It integrates naturally with PyTorch, PIL, NumPy, and supports asynchronous processing, automatic request validation, and real-time interactive documentation via Swagger.

Its design makes it especially suitable for inference services: lightweight, fast to start, and easy to test.

## Use Case: Image Classification with TorchScript

Assume we have trained a convolutional network on Fashion MNIST, exported it as `cnn_scripted.pt`, and want to expose it as a REST API that receives an image and returns its predicted class. Below is the complete implementation:

Listing 21.1: FastAPI service for image inference using TorchScript

```
from fastapi import FastAPI, File, UploadFile
from fastapi.responses import JSONResponse
from PIL import Image
import torch
import torchvision.transforms as transforms
from io import BytesIO

# Device configuration
device = torch.device("cuda" if torch.cuda.is_available()
    else "cpu")
```

```
# Class labels for Fashion MNIST
LABELS = {
    0: "T-shirt/top", 1: "Trouser", 2: "Pullover",
    3: "Dress",        4: "Coat",      5: "Sandal",
    6: "Shirt",        7: "Sneaker",   8: "Bag",      9: "Ankle
        boot"
}

# Load model once on startup
model = torch.jit.load("cnn_scripted.pt", map_location=
    device)
model.eval()

# Define image preprocessing
transform = transforms.Compose([
    transforms.Grayscale(),
    transforms.Resize((28, 28)),
    transforms.ToTensor()
])

app = FastAPI()

@app.post("/predict/")
async def predict(file: UploadFile = File(...)):
    try:
        image = Image.open(BytesIO(await file.read())).convert("RGB")
        x = transform(image).unsqueeze(0).to(device)

        with torch.no_grad():
            output = model(x)
            probs = torch.softmax(output, dim=1)
            pred_idx = probs.argmax(dim=1).item()
            confidence = probs.max(dim=1).item()
            label = LABELS[pred_idx]

    return JSONResponse(content={
```

```
        "prediction_index": pred_idx,
        "label": label,
        "confidence": round(confidence, 4)
    })

except Exception as e:
    return JSONResponse(status_code=400, content={"error": str(e)})
```

This service loads the model once into memory, pre-processes the uploaded image to match the model's input, and returns not only the class index, but also the readable label and its associated confidence score.

## Running the Service Locally

To launch the service locally, save the code to a file named `api.py`, then execute:

Listing 21.2: Starting FastAPI service with Uvicorn

```
uvicorn api:app --host 0.0.0.0 --port 8000 --reload
```

Once running, visit:

```
http://127.0.0.1:8000/docs
```

FastAPI will automatically generate a web-based interface where you can upload images, trigger predictions, and inspect results. This is one of the reasons it is so popular among deep learning practitioners - testing your model becomes as easy as using a website.

## Moving to the Cloud and Beyond

Once your service works locally, you can deploy it to any modern environment. Common options include:

- Hosting on platforms like Render or Hugging Face
- Running in Colab using `ngrok` to expose the local port
- Deploying to a virtual machine on GCP, AWS, or Azure, often within a Docker container

In all cases, the structure remains the same, but the entry point is now a public-facing endpoint ready to serve predictions.

## Stability and Performance Tips

To ensure your API behaves consistently and efficiently in production:

Always set `model.eval()` before inference to ensure deterministic behavior. Wrap the prediction block with `torch.no_grad()` to save memory and disable gradient tracking. Avoid loading the model per request - instead, load it once when the server starts.

Use `map_location=torch.device("cpu")` when loading models on CPU-only environments. Failing to do so can lead to runtime errors if the model was saved on GPU.

For better responsiveness, keep image transforms lightweight and process inputs on the same device as the model.

If you plan to handle multiple requests per second, disable FastAPI's `--reload` flag and use a production WSGI server like Gunicorn or Uvicorn with worker pools. You can also implement batch inference or async task queues using Celery or Ray Serve in more advanced systems.

## Summary

Deploying a deep learning model is not an afterthought - it is the moment the model becomes a system. A simple REST API, powered by

FastAPI and TorchScript, can transform a local script into a reusable, testable, and accessible service.

It doesn't take much code. But it does require clear thinking about interfaces, devices, and stability. Once deployed, your model is no longer just a file - it's part of a product.

*Deploying a model is not just about making it available - it's about turning an idea into a working system.*

## Lesson 22

### Advancing into Specialized Domains

#### Introduction

If you've reached this point, you already have a solid foundation in deep learning. You understand how neural networks are built, trained, optimized, and deployed. You know what it means to serve predictions in real-time, and how to manage models with discipline and care.

But deep learning is not a single field - it is an ecosystem of specialized domains. Each domain comes with its own challenges, model architectures, and conventions. In this lesson, we'll map out the four most common directions for professional specialization: computer vision, natural language processing, tabular learning, and time series analysis.

#### Why Specialization Matters

No single model architecture fits all problems. A model that excels at image recognition may completely fail when applied to text. A method designed for tabular CRM data may struggle to capture temporal dependencies in financial forecasts. Specialization in deep learning is not based on choosing a favorite algorithm - it's about adapting your tools

to the structure of the data and the nature of the task. Understanding this mapping is key to working effectively in applied deep learning.

## Computer Vision

Image data poses unique challenges. Unlike flat vectors, images contain spatial structure, local dependencies, and often large file sizes. Models must extract features hierarchically while preserving spatial relationships.

Early computer vision networks rely on convolutional layers, which scan for patterns in local patches. Architectures like CNNs are foundational for tasks such as image classification, detection, and segmentation. For object localization, architectures like YOLO and DETR introduce mechanisms for bounding box regression and class assignment. For segmentation, models such as U-Net or more recent architectures like SAM provide pixel-level predictions - answering not just what is in the image, but where.

For example, a manufacturing client might want not only to detect if a defect exists, but also to pinpoint its exact location within the image. This isn't classification - it's structured output over space, and requires a specialized model.

## Language Models

Text data is sequential, context-dependent, and semantically rich. Unlike images, it cannot be interpreted by pixel structure but instead depends on token relationships, syntax, and meaning. The core architecture that revolutionized NLP is the Transformer - a model capable of capturing long-range dependencies and parallelizing sequence modeling.

Transformers serve as the backbone of major language models such as BERT, RoBERTa, T5, GPT, and LLaMA. These models are pre-trained on massive corpora and then fine-tuned for downstream tasks like classification, summarization, or dialogue.

In some applications, information retrieval becomes essential. In such cases, architectures that support RAG (Retrieval-Augmented Generation) are used - combining a language model with a search component. In more interactive scenarios, agent-based systems allow the model not only to generate text but also to take actions or query external tools.

A typical use case would be a customer service system that extracts intent from a support email or generates a concise summary of a legal document. The model needs to understand, condense, and act - all based on unstructured textual input.

## Tabular Data

Tabular data is ubiquitous in business: sales logs, CRM exports, finance records, and customer metadata. While the format seems simple - rows and columns - the modeling challenges are real. Categorical features, missing values, imbalanced classes, and high-cardinality identifiers complicate the learning process.

Traditional ensemble models like XGBoost and LightGBM still dominate many tabular competitions. They are fast, robust, and interpretable. However, neural approaches are evolving. MLPs work well when data volume is high, and architectures like TabNet or FT-Transformer introduce attention mechanisms to model column-wise relationships more effectively.

In domains like credit scoring or churn prediction, explainability is not optional. Regulatory and operational constraints often require you

to justify why a prediction was made - for example, using SHAP values or permutation-based methods.

## Time Series and Sequential Forecasting

Temporal data introduces a new axis: time. Each prediction depends not just on features, but on what came before. Time series forecasting often involves trends, seasonality, delays, or missing intervals - requiring models to learn both sequence and context.

For this, we use recurrent models like LSTM and GRU, which maintain internal memory over time. CNN-based models like TCN (Temporal Convolutional Networks) filter over time windows, providing an alternative to recurrence. More recently, Transformer variants such as TFT or Autoformer offer scalable alternatives for long-range sequence modeling. For high-frequency signals, like sensor data or control systems, State Space Models - such as Mamba - are emerging as powerful alternatives.

Imagine forecasting product demand across hundreds of SKUs, each with its own seasonal pattern, promotional history, and geographic variation. The goal is not just to remember the past - but to learn how it evolves and when it shifts.

## How to Choose a Path Forward

Each domain comes with distinct considerations - not just in data structure, but in modeling goals, inference constraints, and system integration. There is no one-size-fits-all.

A simple rule of thumb:

- If your input is an image, start with a convolutional model.
- If you're working with free-form text, look into Transformers.

- If you have structured rows and columns, consider ensemble models or MLPs.

- If your data evolves over time, explore sequential architectures - recurrent or transformer-based.

As always, the structure of the data should guide your choice, not the popularity of the algorithm.

## Summary

Deep learning is a family of approaches that are tailored to different types of data and different modeling challenges. Vision, language, tables, and sequences - each requires its own way of thinking.

Choosing a direction for specialization depends on your domain, your goals, and the nature of the problem at hand. Some engineers become experts in one modality. Others learn to blend them.

*Choose your tools by the shape of your data, not the fashion of the field.*

## Lesson 23

### Roadmap for the Industrial DL Engineer

Congratulations - you've reached the final lesson of the foundational deep learning course. We've come a long way together: from the inner workings of a neural network to deploying a real-time prediction service in the cloud.

Now it's time to pause and ask a deeper question: what does the path forward look like for a deep learning engineer in the real world?

This lesson is not about new code or new models. It's about mindset. It's about building systems that last, making decisions you can stand behind, and taking your work from experiment to impact.

### From Model to Production System

Training an accurate model is only the beginning. In industrial settings, accuracy alone is not enough. A good model is one that can be deployed, monitored, updated, and trusted over time.

This means shifting your mindset from “the model” to “the system.” A file with weights is not a solution. A full pipeline - from training and validation to service and monitoring - is what makes deep learning truly operational.

## Three Core Components of a DL System

Every real-world deep learning system has three foundational components:

The first is the training component. This includes data pipelines, hyperparameter configuration, model architecture, and reproducibility. A well-designed training pipeline allows experiments to be repeated months later - with the same results.

The second is the validation and monitoring component. This is what tells you whether the model still works on new data, whether performance is drifting, and whether action is needed. Without ongoing monitoring, a once-good model can become silently harmful.

The third is the interface - usually a REST API or streaming service - through which predictions are consumed. This is what other systems interact with. It should be fast, stable, and transparent.

These three parts must be aligned, versioned, and auditable. Together, they form the foundation of any production ML system.

## Case Study: Email Intent Classification

Imagine you're tasked with building a model that detects user intent from incoming customer emails.

- A training module that logs every experiment, including:
  - Dataset version
  - Hyperparameters
  - Final metrics
- A validation suite that:
  - Tests the model weekly on fresh support tickets

- Alerts on drops in accuracy or coverage
- An API implemented that returns:
  - Predictions
  - Confidence scores
  - Short textual justifications

All of this would be tracked with MLflow or a custom Git-based metadata system, ensuring that every change is documented and reversible. This is the standard of production-grade deep learning.

## Best Practices for Real-World Engineers

To work professionally in deep learning, it is not enough to write good code. You need discipline. Always document what you do. Every model version should have a clear date, configuration, and evaluation results. Ensure reproducibility. Someone, often you, would be able to rerun the same training job in six months and get the same output. Use robust serialization formats like TorchScript, ONNX, or SavedModel. Avoid custom scripts or ad-hoc saving.

Never skip validation. A model that performs inconsistently, or that changes behavior over time, cannot be trusted. Design modular pipelines. Your API should not care if the underlying model changes. Your training logic should be isolated from serving logic. And always evaluate for consistency over time, not just for single-point accuracy.

## Ethics and Responsibility

Deep learning can influence hiring, diagnosis, pricing, and more. With power comes responsibility. Never deploy a model you haven't tested

across multiple populations. Always understand where your data came from - and what it may be missing. Be prepared to explain predictions, even if only partially. And most importantly: ask yourself whether you would trust this system if it affected you or someone you care about.

Responsible AI is not a legal checkbox. It's an engineering value.

## True Professionalism

What separates an engineer from a hobbyist is not only the depth of their models - but the depth of their thinking. A professional deep learning engineer doesn't just tune hyperparameters. They build systems that are resilient, reproducible, observable, and safe. It takes attention to detail. It takes humility. But it's also what transforms a promising model into a real solution.

## DL Engineering Manifesto

The following principles summarize the mindset and values of a professional deep learning engineer operating in the real world. These are not rules-they are field-tested habits that distinguish mature systems from fragile prototypes.

- **Every model is a system.** A model is only as good as the data pipeline, validation process, and deployment environment it lives in.
- **Reproducibility over intuition.** Every experiment should be traceable. What cannot be repeated cannot be trusted.
- **Version everything.** Data, code, models, configuration, metrics track them all. Today's shortcut is tomorrow's mystery bug.

- **Validate beyond accuracy.** Evaluate drift, stability, robustness, and confidence-not just a single metric on a static test set.
- **Design for failure.** Monitoring is not optional. Assume things will break. Build alerts, retries, and fallback strategies.
- **Explainability matters.** If a model makes a decision you can't explain, it shouldn't be making real-world decisions.
- **Code is not the product-impact is.** The goal isn't to write fancy models; it's to solve real problems with dependable tools.
- **Choose tools by data, not trend.** Let the shape of the problem guide your architecture-not what's popular on arXiv.
- **Bias is a bug.** Actively search for it. Mitigate it early. Document your mitigation.
- **Build trust.** Your job isn't just to optimize metrics-it's to create systems that others can rely on.

These are the principles that will guide you as you grow from building models to engineering reliable systems. Keep them close.

## Where to Go From Here

This course gave you the foundation. You now have the vocabulary, the tools, and the confidence to start real projects.

From here, you can go deeper: into computer vision, language models, tabular learning, or time series. You can specialize, or you can integrate. You can work solo, or join a team. But whatever path you take, you now carry the mindset that matters: methodical, ethical, and applied.

## Summary

A career in deep learning evolves constantly. But the principles of sound engineering remain: learn, build, measure, share. Choose a small, real project. Treat it seriously. Track your assumptions. Deploy your ideas. I hope this course gave you not only knowledge - but the confidence, tools, and motivation to continue. The world needs the systems you'll build - and you're ready.

Good Luck,

Barak