

Neural Networks from Math

Vu Hung Nguyen

2025-11-04

Contents

Preface	iii
1 Overview	1
2 Functions and Graphs	3
2.1 From Tables to Rules	3
2.2 Exercises	4
3 Data and Error	5
3.1 Measuring Error	5
3.2 Exercises	6
4 Vectors and Matrices	7
4.1 Notation	7
4.2 Exercises	8
5 Linear Models	9
5.1 Form and Geometry	9
5.2 Exercises	10
6 Nonlinearity and Activations	11
6.1 Common Activations	11
6.2 Kernel Functions	12
6.3 Exercises	12
7 Composing Layers	13
7.1 Two-Layer Example	13
7.2 Feature Learning Functions	14
7.3 Exercises	15
8 Loss and Optimisation	16
8.1 Objective Functions	16
8.2 Loss Choices	17
8.3 Cost Functions	18
8.4 Regularisation Functions	18
8.5 Gradient Descent (Conceptually)	19
9 Training Loop Concepts	21
9.1 The Loop	21
9.2 Exercises	22

10 A Simple Neural Network Architecture	23
10.1 Shapes and Flow	23
10.2 Exercises	24
11 Limits and Ethics	25
11.1 Limits	25
11.2 Ethics	25
11.3 Further Directions	26
12 Glossary	27
A Worksheets: Checkpoints	28
B Worksheets: Exercises	29
References	30
About the Author	32

Preface

This brief book takes an example-first path from familiar school mathematics to the core ideas behind a simple neural network architecture. It is designed for motivated high-school students: short chapters, visual intuition, and concrete checkpoints.

Chapter 1

Overview

This book takes you from familiar mathematics to a working understanding of simple neural networks. It is written in clear British English, and designed for undergraduates and motivated high-school students.

Learning Objectives

1. Understand the book’s roadmap and how chapters connect.
2. Recognise the key ingredients of a neural network.
3. Learn how examples, metaphors, and visuals support intuition.

Roadmap

We move from functions and graphs to data and error, then to vectors and matrices that compactly represent many inputs. Linear models give us a first predictive rule; nonlinear activations add the “kinks” that let us model more interesting patterns. Stacking layers composes these ideas. Finally, we discuss loss and optimisation, the training loop, a simple network, and practical limits and ethics. The foundations we build connect to the rich history of neural networks [5, 8] and the modern deep learning revival [2].

Example 1.1

Consider predicting house price from size. A function maps size (input) to price (output). A linear model draws a straight line; a neural network allows bends via activations and layers, improving fit when reality is not a straight line.

Remark 1.1

Throughout, we prioritise intuition first, then formalism. Visuals accompany core definitions where helpful.

How to Use This Book

Skim the learning objectives at the start of each chapter. Work through examples; attempt exercises before revealing the upside-down hints. Use the glossary to refresh key terms.

Chapter 2

Functions and Graphs

Functions describe how inputs map to outputs. Graphs help us see this mapping.

Learning Objectives

1. Interpret a function as a rule from input to output.
2. Read tables and plots to understand trends.
3. Distinguish linear from nonlinear patterns visually.

2.1 From Tables to Rules

Suppose we record study time (hours) and score (percentage). A function f turns an input x into an output $y = f(x)$. A straight trend suggests a linear rule; bends suggest nonlinearity.

Definition 2.1

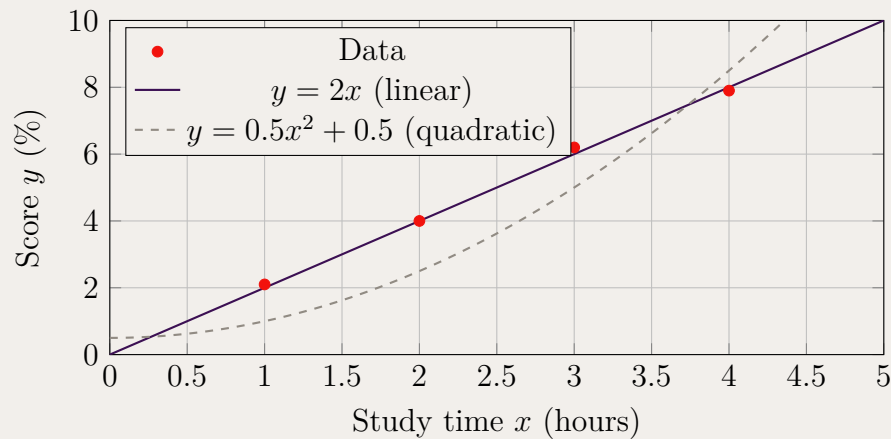
A *function* is a mapping assigning each input a single output. We often visualise $(x, f(x))$ as points on a graph.

Example 2.1

Measured study time (hours) and score (%):

x	1	2	3	4
y	2.1	4.0	6.2	7.9

If doubling x roughly doubles y , a straight line is sensible. Below we plot the points, a simple linear rule $y = 2x$, and a curved (quadratic) rule $y = 0.5x^2 + 0.5$ to contrast linear vs nonlinear behaviour.



The linear rule tracks the trend closely; the quadratic bends upward and will diverge from the roughly proportional pattern as x grows.

2.2 Exercises

Given points $(1, 2)$, $(2, 4.1)$, $(3, 5.9)$, would a straight line be a reasonable model? Explain briefly.

Hint: Turn the page upside down to read Yes. The ratios are consistent with a near-linear trend; noise causes small deviations.

Chapter 3

Data and Error

Predictions meet reality through data. Error tells us how far off we are.

Learning Objectives

1. Define prediction, target, and error for a dataset.
2. Compute simple average absolute and squared error.
3. Explain why averaging error stabilises noisy measurements.

3.1 Measuring Error

Given inputs x_i with observed outputs y_i and predictions \hat{y}_i , an error for item i is $e_i = \hat{y}_i - y_i$. Two popular summaries are the mean absolute error (MAE) and mean squared error (MSE):

$$\text{MAE} = \frac{1}{n} \sum_i |e_i|, \quad \text{MSE} = \frac{1}{n} \sum_i e_i^2.$$

Example 3.1

Using the study–score data from the previous chapter

x	1	2	3	4
y	2.1	4.0	6.2	7.9

consider the linear rule $\hat{y} = 2x$. The predictions and errors are:

x	1	2	3	4
$\hat{y} = 2x$	2.0	4.0	6.0	8.0
$e = \hat{y} - y$	0.1	0.0	-0.2	0.1
$ e $	0.1	0.0	0.2	0.1
e^2	0.01	0.00	0.04	0.01

Thus

$$\text{MAE} = \frac{0.1+0.0+0.2+0.1}{4} = 0.1, \quad \text{MSE} = \frac{0.01+0.00+0.04+0.01}{4} = 0.015.$$

The small MAE/MSE values indicate the line $\hat{y} = 2x$ matches the trend closely for these points.

Remark 3.1

Squaring emphasises large mistakes; absolute value treats all deviations proportionally. Choice depends on the application.

3.2 Exercises

For true values $y = (1, 2, 3)$ and predictions $\hat{y} = (1.2, 1.9, 3.4)$, compute MAE and MSE.

Hint: Errors: $(0.2, -0.1, 0.4)$. $\text{MAE} = (0.2 + 0.1 + 0.4)/3 = 0.233\dots$; $\text{MSE} = (0.04 + 0.01 + 0.16)/3 = 0.07\dots$

Chapter 4

Vectors and Matrices

Vectors collect features; matrices collect many vectors and define linear transformations.

Learning Objectives

1. Represent data points as vectors of features.
2. Use matrix–vector multiplication to combine features linearly.
3. Interpret a matrix as a transformation of space.

4.1 Notation

Write a feature vector as $\mathbf{x} = (x_1, \dots, x_d)^\top$. A weight vector \mathbf{w} and bias b define a linear score $s = \mathbf{w}^\top \mathbf{x} + b$.

Definition 4.1

A *matrix* $W \in \mathbb{R}^{m \times d}$ applied to a vector $\mathbf{x} \in \mathbb{R}^d$ yields $W\mathbf{x} \in \mathbb{R}^m$, combining inputs into m outputs.

Example 4.1

With two features (height, width) and two outputs (sum, difference), $W = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ maps a vector $\mathbf{x} = (x_1, x_2)^\top$ to $W\mathbf{x} = (x_1 + x_2, x_1 - x_2)^\top$. For a concrete calculation, let $\mathbf{x} = (3, 2)^\top$. Then

$$W\mathbf{x} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \cdot 3 + 1 \cdot 2 \\ 1 \cdot 3 + (-1) \cdot 2 \end{bmatrix} = \begin{bmatrix} 5 \\ 1 \end{bmatrix}.$$

So the first output (sum) is 5 and the second output (difference) is 1.

4.2 Exercises

If $W = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$ and $\mathbf{x} = (1, 2)^\top$, compute $W\mathbf{x}$.

Hint: $(2, 6)^\top$.

Chapter 5

Linear Models

A linear model scores inputs by a weighted sum plus a bias; decisions follow from thresholds or regression.

Learning Objectives

1. Write a linear model as $\hat{y} = \mathbf{w}^\top \mathbf{x} + b$.
2. Explain the geometric view: lines and hyperplanes.
3. Recognise when linear models are insufficient.

5.1 Form and Geometry

The set of points with constant score forms a hyperplane. When classes are linearly separable, a single hyperplane can divide them.

Example 5.1

Classify points above a line: take $\hat{y} = 2x_1 - x_2 + 0.5$ and predict positive when $\hat{y} > 0$.

For $(x_1, x_2) = (2, 1)$,

$$\hat{y} = 2 \cdot 2 - 1 + 0.5 = 3.5 > 0 \Rightarrow \text{positive.}$$

For $(x_1, x_2) = (0, 1)$,

$$\hat{y} = 2 \cdot 0 - 1 + 0.5 = -0.5 < 0 \Rightarrow \text{negative.}$$

The decision boundary $\hat{y} = 0$ is the straight line $x_2 = 2x_1 + 0.5$.

Remark 5.1

When data is not linearly separable, linear models struggle. One alternative approach uses *kernel functions* (e.g., in Support Vector Machines) to implicitly map inputs to higher-dimensional spaces where separation becomes easier. Neural networks handle nonlinearity differently, through explicit nonlinear activations (see Chapter 5).

5.2 Exercises

Find a line that separates points A: $(0, 2), (1, 3)$ from B: $(2, 0), (3, 1)$ if possible.

Hint: Try $x_2 = x_1 + 1.5$. Points A lie above, B below.

Chapter 6

Nonlinearity and Activations

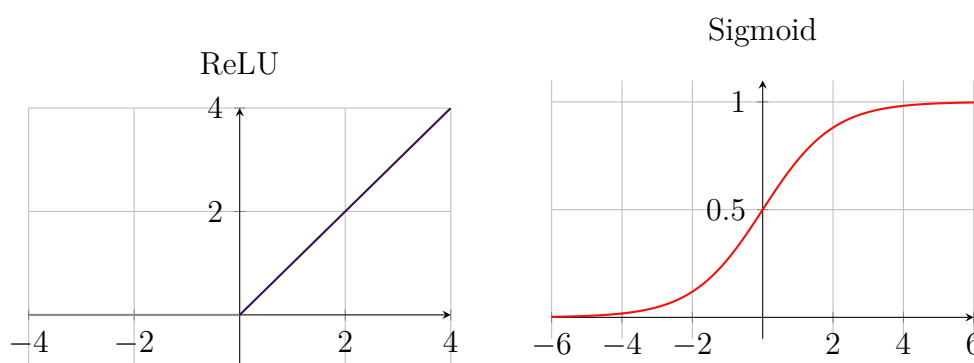
Nonlinear activations introduce bends that let models capture curved patterns.

Learning Objectives

1. Describe common activations: ReLU, sigmoid, tanh.
2. Explain why stacking linear layers without activations stays linear.
3. Match activations to tasks (e.g., sigmoid for probabilities).
4. Understand kernel functions as an alternative approach to nonlinearity.

6.1 Common Activations

ReLU: $\text{ReLU}(z) = \max(0, z)$ adds a kink at zero [6]. Sigmoid $\sigma(z) = 1/(1 + e^{-z})$ squashes to $(0, 1)$. Tanh squashes to $(-1, 1)$.



Why this matters. Activations create nonlinearity so networks can approximate curved relationships. ReLU is simple and keeps gradients flowing for positive inputs, aiding optimisation. Sigmoid maps real numbers to $(0, 1)$, ideal for probabilities; but it can saturate (flat tails), slowing learning when inputs are very large in magnitude.

Bio example (sigmoid). The logistic (sigmoid) curve models population growth with carrying capacity and neuron firing rates as a function of membrane potential: response is low for small input, steep near threshold, and saturates at high input.

Remark 6.1

Linear \circ linear is still linear. Nonlinearity between linear steps is essential.

6.2 Kernel Functions

Kernel functions offer an alternative approach to handling nonlinearity. The *kernel trick* implicitly maps input data to higher-dimensional feature spaces, where linear separation becomes possible without explicitly computing the transformed features.

Definition 6.1

A *kernel function* $k(\mathbf{x}, \mathbf{x}')$ computes the inner product of feature vectors in a transformed space: $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$, where ϕ maps inputs to a higher-dimensional space. Common kernels include the polynomial kernel and the radial basis function (RBF) kernel.

Example 6.1

For data that is not linearly separable in 2D (e.g., points arranged in concentric circles), a kernel can map them to 3D where a plane can separate them. The kernel trick allows this transformation without explicitly computing the 3D coordinates, making it computationally efficient.

Remark 6.2

Neural networks use *explicit activations* instead of kernels. Unlike kernels, which provide a fixed transformation, activations are learned during training—each layer's weights and activations adapt to the data. This flexibility allows networks to discover task-specific nonlinear transformations rather than relying on predefined kernel functions. Additionally, activations enable deep, composable architectures that can learn hierarchical features.

6.3 Exercises

Sketch and compare $\tanh(z)$ and leaky ReLU $\text{LReLU}_\alpha(z) = \max(\alpha z, z)$ with $\alpha = 0.1$. Where are they most/least sensitive, and what ranges do their outputs take?

Hint: ReLU is flat for $z < 0$ and slope 1 for $z > 0$. Sigmoid is steepest near 0, flat in tails.

Chapter 7

Composing Layers

Stack linear steps with activations to form complex functions.

Learning Objectives

1. Define a layer as $\mathbf{h} = g(W\mathbf{x} + \mathbf{b})$.
2. Explain how two layers can form piecewise linear curves.
3. Understand hidden units as learned features.
4. Describe feature learning functions like PCA and autoencoders.

7.1 Two-Layer Example

Let $\mathbf{h} = \text{ReLU}(W_1\mathbf{x} + \mathbf{b}_1)$ and $\hat{y} = W_2\mathbf{h} + b_2$. Even with ReLU, the output is a flexible piecewise linear function of \mathbf{x} .

Example 7.1

In one dimension, two ReLUs can create a “bump” by combining kinks at different locations. Consider a tiny network with one input x , two hidden units, and one output:

$$\mathbf{h} = \text{ReLU}(W_1 x + \mathbf{b}_1), \quad \hat{y} = W_2 \mathbf{h} + b_2,$$

where

$$W_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad \mathbf{b}_1 = \begin{bmatrix} -1 \\ -3 \end{bmatrix}, \quad W_2 = \begin{bmatrix} 1 & -1 \end{bmatrix}, \quad b_2 = 0.$$

Thus the hidden units are

$$h_1 = \text{ReLU}(x - 1), \quad h_2 = \text{ReLU}(-x - 3) = \text{ReLU}(-(x + 3)).$$

The output becomes $\hat{y} = h_1 - h_2$. Evaluate at a few inputs:

x	$h_1 = \max(0, x - 1)$	$h_2 = \max(0, -x - 3)$	$\hat{y} = h_1 - h_2$
-5	0	2	-2
-3	0	0	0
0	0	0	0
2	1	0	1
4	3	0	3

Interpretation:

- For $x \leq -3$, only the second unit is active, so $\hat{y} = -h_2$ decreases linearly.
- For $-3 < x < 1$, both units are off, so $\hat{y} = 0$.
- For $x \geq 1$, only the first unit is active, so $\hat{y} = h_1$ increases linearly.

This piecewise linear shape forms a flat region with rising and falling flanks—one way to build a “bump” by composing ReLUs.

7.2 Feature Learning Functions

Feature learning functions extract useful representations from data. Hidden layers in neural networks learn features automatically; simpler methods like PCA provide linear transformations for dimensionality reduction.

Definition 7.1

A *feature learning function* transforms raw inputs into a representation that captures essential patterns, making subsequent tasks (classification, prediction) easier. The learned features can reduce dimensionality or highlight important relationships in the data.

Example 7.2: Principal Component Analysis (PCA)

PCA finds directions (principal components) that capture maximum variance. Given data points $\mathbf{x}_1 = (1, 2)$, $\mathbf{x}_2 = (2, 3)$, $\mathbf{x}_3 = (3, 4)$, first centre them: $\bar{\mathbf{x}} = (2, 3)$. Centred points are $(-1, -1)$, $(0, 0)$, $(1, 1)$. The first principal component aligns with the direction $(1, 1)$, capturing the main variation. Projecting onto this component reduces 2D data to 1D while preserving the most information.

Example 7.3: Autoencoders

An autoencoder is a neural network trained to reconstruct its input. It has an encoder (maps input \mathbf{x} to hidden representation \mathbf{h}) and a decoder (maps \mathbf{h} back to $\hat{\mathbf{x}}$). By constraining the hidden layer to be smaller than the input, the network learns a compressed representation. For instance, with 10 inputs, a hidden layer of 3 units forces the network to learn a 3D summary of the essential information needed for reconstruction.

Remark 7.1

PCA provides fixed linear transformations; autoencoders learn nonlinear feature extraction through training. Hidden layers in standard neural networks also learn features, but they are optimised for the end task (e.g., classification) rather than reconstruction. A foundational result shows that multilayer feedforward networks are universal approximators [3]—they can approximate any continuous function given enough hidden units.

If PCA reduces 100D data to 10D, what percentage of the original dimensions are retained? What does this suggest about the data?

Hint: 10% of dimensions retained. If PCA works well, it suggests the data has strong correlations and can be compressed without much loss.

7.3 Exercises

How many linear regions can a sum of two shifted ReLUs create on the line?

Hint: Up to three regions separated by the two kink locations.

Chapter 8

Loss and Optimisation

Loss quantifies mismatch between predictions and targets; optimisation reduces loss.

Remark 8.1

The terms ‘loss function’ and ‘cost function’ are used interchangeably; both measure how far predictions are from targets. Both are types of *objective functions*—what we optimise during training. This book uses ‘loss’ throughout for consistency.

Learning Objectives

1. Understand objective functions as the general goal of optimisation.
2. Define a loss function suitable for a task.
3. Explain how regularisation prevents overfitting.
4. Describe gradient descent at a high level.
5. Relate learning rate to step size and stability.

8.1 Objective Functions

An objective function represents the goal of the learning algorithm—what we aim to optimise (minimise or maximise) during training. Loss functions and cost functions are specific types of objective functions used in machine learning.

Definition 8.1

An *objective function* is a mathematical expression that quantifies how well a model performs. In neural networks, we typically minimise the objective function to find optimal parameters θ that best fit the training data.

Remark 8.2

The hierarchy is: *objective function* (general term) encompasses *loss function* (per example) and *cost function* (aggregated over dataset). In practice, ‘loss’ and ‘cost’ are often used interchangeably, and both are types of objective functions we minimise.

Example 8.1

For a regression task, the objective might be to minimise prediction error. The mean squared error (MSE) serves as both a loss function (for individual predictions) and a cost function (when averaged over all examples). The learning algorithm’s goal is to find parameters that minimise this objective.

8.2 Loss Choices

For regression, MSE is common; for binary classification, logistic loss pairs well with sigmoid outputs.

Example 8.2: Regression loss (MSE)

Suppose a model predicts house prices (in thousands) for three homes: $\hat{y} = (210, 195, 250)$, and true prices are $y = (200, 205, 240)$. Errors are $e = \hat{y} - y = (10, -10, 10)$. The mean squared error is

$$\text{MSE} = \frac{1}{3}(10^2 + (-10)^2 + 10^2) = \frac{1}{3}(100 + 100 + 100) = 100.$$

Units are squared (here, thousands²). A smaller value indicates closer predictions.

Example 8.3: Logistic loss

For a binary label $y \in \{0, 1\}$ with sigmoid output $p = \sigma(z)$ interpreted as $\Pr(y = 1 \mid x)$, the logistic loss for a single example is

$$\ell(y, p) = -(y \log p + (1 - y) \log(1 - p)).$$

This loss function is rooted in logistic regression [1]. Consider two cases with $p = 0.9$:

- If $y = 1$: $\ell = -(1 \cdot \log 0.9 + 0 \cdot \log 0.1) \approx 0.105$ (small penalty).
- If $y = 0$: $\ell = -(0 \cdot \log 0.9 + 1 \cdot \log 0.1) \approx 2.303$ (large penalty for confident wrong prediction).

This asymmetry encourages calibrated probabilities: confident and correct is rewarded; confident and wrong is penalised heavily.

8.3 Cost Functions

A cost function aggregates the loss over the entire dataset, usually by averaging the individual losses. It is minimised during training to find optimal model parameters.

Definition 8.2

A *cost function* $J(\boldsymbol{\theta})$ computes the average loss across all training examples:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i(\boldsymbol{\theta})),$$

where ℓ is the loss function for a single example, n is the number of examples, and $\boldsymbol{\theta}$ represents the model parameters.

Example 8.4

Using the same house price data from earlier, suppose we have predictions $\hat{y} = (210, 195, 250)$ and true values $y = (200, 205, 240)$ for three homes. The cost function (mean squared error) aggregates the individual squared errors:

$$J = \frac{1}{3} \sum_{i=1}^3 (y_i - \hat{y}_i)^2 = \frac{1}{3} (100 + 100 + 100) = 100.$$

During training, we adjust parameters to minimise this cost; a lower cost indicates better overall fit to the data.

8.4 Regularisation Functions

Regularisation functions penalise model complexity to prevent overfitting. They are added to the cost function, encouraging the model to favour simpler parameter values.

Definition 8.3

A *regularised cost function* combines the original cost with a regularisation term:

$$J_{\text{reg}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda \cdot R(\boldsymbol{\theta}),$$

where $J(\boldsymbol{\theta})$ is the original cost, $\lambda > 0$ is the regularisation strength (hyperparameter), and $R(\boldsymbol{\theta})$ is the regularisation function that penalises complexity.

Example 8.5: L2 regularisation (Ridge)

L2 regularisation penalises the sum of squared parameters: $R_{L2}(\boldsymbol{\theta}) = \sum_j \theta_j^2$. For parameters $\boldsymbol{\theta} = (2, -3, 1.5)$ with $\lambda = 0.1$, the regularisation term is

$$\lambda \cdot R_{L2} = 0.1 \cdot (2^2 + (-3)^2 + 1.5^2) = 0.1 \cdot (4 + 9 + 2.25) = 1.525.$$

If the original cost is $J = 50$, the regularised cost becomes $J_{\text{reg}} = 50 + 1.525 = 51.525$. Larger parameters increase the penalty, encouraging smaller values.

Example 8.6: L1 regularisation (Lasso)

L1 regularisation penalises the sum of absolute parameter values: $R_{L1}(\boldsymbol{\theta}) = \sum_j |\theta_j|$. For the same parameters $\boldsymbol{\theta} = (2, -3, 1.5)$ with $\lambda = 0.1$:

$$\lambda \cdot R_{L1} = 0.1 \cdot (|2| + |-3| + |1.5|) = 0.1 \cdot (2 + 3 + 1.5) = 0.65.$$

With $J = 50$, the regularised cost is $J_{\text{reg}} = 50 + 0.65 = 50.65$. L1 tends to drive some parameters exactly to zero, promoting sparsity.

Remark 8.3

L2 regularisation (Ridge) shrinks parameters smoothly; L1 regularisation (Lasso) can zero out parameters entirely. Both help prevent overfitting by discouraging overly large parameter values. Dropout [10] is another regularisation technique that randomly sets some neurons to zero during training, forcing the network to learn more robust features.

Given $\boldsymbol{\theta} = (1, -2, 0.5)$ and $\lambda = 0.2$, compute both L1 and L2 regularisation terms. Which is larger?

Hint: L2: $0.2 \cdot (1^2 + (-2)^2 + 0.5^2) = 0.2 \cdot (1 + 4 + 0.25) = 1.05$. L1: $0.2 \cdot (1 + 2 + 0.5) = 0.7$. L2 is larger because squaring amplifies larger values.

8.5 Gradient Descent (Conceptually)

Imagine standing on a landscape where height is loss. At each step, move a little in the downhill direction; repeat until you are low enough. The theoretical foundation dates to stochastic approximation methods [7].

Example 8.7: One variable

Let $f(x) = (x - 3)^2$. Its derivative is $f'(x) = 2(x - 3)$. Starting at $x_0 = 0$ with learning rate $\eta = 0.5$, gradient descent updates are

$$x_{k+1} = x_k - \eta f'(x_k) = x_k - 0.5 \cdot 2(x_k - 3) = 3 - (x_k - 3).$$

Numerically: $x_1 = 1.5$, $x_2 = 2.25$, $x_3 = 2.625 \dots$ which approaches the minimiser $x^* = 3$.

Example 8.8: Two variables

Let $f(x, y) = (x - 1)^2 + (y + 2)^2$. The gradient is $\nabla f = (2(x - 1), 2(y + 2))$. With $(x_0, y_0) = (0, 0)$ and $\eta = 0.25$:

$$(x_1, y_1) = (0, 0) - 0.25 (2(-1), 2(2)) = (0.5, -1).$$

Next step uses the new gradient: $\nabla f(x_1, y_1) = (2(-0.5), 2(1)) = (-1, 2)$, thus

$$(x_2, y_2) = (0.5, -1) - 0.25 (-1, 2) = (0.75, -1.5).$$

The iterates head toward the minimiser $(1, -2)$.

Many inputs (neural nets). For networks, parameters form a long vector θ (all weights and biases). The loss $L(\theta)$ is averaged over a mini-batch. We compute the gradient $\nabla L(\theta)$ efficiently by backpropagation [9] and update

$$\theta_{k+1} = \theta_k - \eta \nabla L(\theta_k).$$

Different layers receive different components of the same gradient, moving all parameters a little in the direction that most reduces loss on the current batch. Modern optimisers like Adam [4] adapt the learning rate per parameter, often converging faster than basic gradient descent.

If steps are too large, what failure can occur?

Hint: You can overshoot and bounce around (diverge) instead of settling.

Chapter 9

Training Loop Concepts

Predict, measure, adjust, repeat. Split data to measure progress fairly.

Learning Objectives

1. Describe the predict–loss–update loop.
2. Explain train/validation/test splits.
3. Recognise overfitting and how validation helps.

9.1 The Loop

For each mini-batch: compute predictions, compute loss, update parameters a little. After each epoch, check validation loss.

Example 9.1

Consider a simple linear model $\hat{y} = wx + b$ trained with mean squared error on one mini-batch of two points: $(x, y) \in \{(1, 2), (2, 4)\}$. Start from $w_0 = 1.0$, $b_0 = 0.0$, learning rate $\eta = 0.1$.

Predictions. With $(w, b) = (1, 0)$, predictions are $\hat{y} = (1, 2)$. Errors $e = \hat{y} - y = (-1, -2)$.

Loss. $\text{MSE} = \frac{1}{2}((-1)^2 + (-2)^2) = \frac{1}{2}(1 + 4) = 2.5$.

Gradients. For MSE with this batch,

$$\frac{\partial L}{\partial w} = \frac{1}{2} \sum 2e_i x_i = e_1 x_1 + e_2 x_2 = (-1) \cdot 1 + (-2) \cdot 2 = -5,$$

$$\frac{\partial L}{\partial b} = \frac{1}{2} \sum 2e_i = e_1 + e_2 = -3.$$

Update. Gradient descent gives

$$w_1 = w_0 - \eta \frac{\partial L}{\partial w} = 1 - 0.1 \cdot (-5) = 1.5, \quad b_1 = b_0 - \eta \frac{\partial L}{\partial b} = 0 - 0.1 \cdot (-3) = 0.3.$$

Repeating on the next mini-batch (or another epoch) will continue to reduce loss; validation loss is checked after full-dataset passes to monitor overfitting.

Remark 9.1

When validation loss rises while training loss falls, you may be overfitting.

9.2 Exercises

Why is a separate test set needed if validation loss looks good?

Hint: To estimate performance on truly unseen data and avoid tuning to the validation set.

Chapter 10

A Simple Neural Network Architecture

A minimal feed-forward network: inputs, a hidden layer, and an output. The earliest computational models of neural activity [5] and the perceptron [8] laid groundwork for these architectures.

Learning Objectives

1. Sketch a tiny network with one hidden layer.
2. Track shapes through layers.
3. Explain forward pass at a high level.

10.1 Shapes and Flow

With d inputs, h hidden units, and one output: $W_1 \in \mathbb{R}^{h \times d}$, $\mathbf{b}_1 \in \mathbb{R}^h$, activation g ; then $W_2 \in \mathbb{R}^{1 \times h}$, $b_2 \in \mathbb{R}$.

Example 10.1

Let $d = 2, h = 3$. Then W_1 has 3 rows and 2 columns; W_2 has 1 row and 3 columns. Compute $\mathbf{h} = g(W_1\mathbf{x} + \mathbf{b}_1)$, then $\hat{y} = W_2\mathbf{h} + b_2$.

Take a numerical instance with ReLU activation $g = \text{ReLU}$:

$$\mathbf{x} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad W_1 = \begin{bmatrix} 1 & -1 \\ 0.5 & 1 \\ 2 & 0 \end{bmatrix}, \quad \mathbf{b}_1 = \begin{bmatrix} 0 \\ -1 \\ 0.5 \end{bmatrix}.$$

First layer pre-activation and activation:

$$\mathbf{z}_1 = W_1\mathbf{x} + \mathbf{b}_1 = \begin{bmatrix} 1 & -1 \\ 0.5 & 1 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 4.5 \end{bmatrix}, \quad \mathbf{h} = g(\mathbf{z}_1) = \text{ReLU}(\mathbf{z}_1) = \begin{bmatrix} 1 \\ 1 \\ 4.5 \end{bmatrix}.$$

Output layer parameters and prediction:

$$W_2 = \begin{bmatrix} 1 & -2 & 0.5 \end{bmatrix}, \quad b_2 = 0.2, \quad \hat{y} = W_2\mathbf{h} + b_2 = 1 \cdot 1 + (-2) \cdot 1 + 0.5 \cdot 4.5 + 0.2 = 1.45.$$

Thus with this small network and input $\mathbf{x} = (2, 1)^\top$, the hidden representation is $\mathbf{h} = (1, 1, 4.5)^\top$ and the scalar output is $\hat{y} = 1.45$.

10.2 Exercises

If $g = \text{ReLU}$ and $\mathbf{h} = (0, 2, 0)^\top$, what is \hat{y} when $W_2 = (1, 1, 1)$ and $b_2 = 0$?

Hint: 2.

Chapter 11

Limits and Ethics

Models have limits. Responsible practice matters.

Learning Objectives

1. Identify common failure modes: overfitting, data shift, bias.
2. Describe simple mitigations students can apply.
3. Reflect on responsible communication of model results.

11.1 Limits

Overfitting memorises training noise; distribution shift breaks assumptions; biased data leads to unfair outcomes.

Example 11.1: Overfitting

Fit a 9th-degree polynomial to 10 noisy points sampled from the straight line $y = 2x$. On the training points, the curve can wiggle to achieve near-zero error; on fresh test points from the same line, error is high because the wiggles chase noise rather than signal.

Example 11.2: Underfitting

Fit a straight line to data generated by a clear U-shaped quadratic like $y = x^2$ with small noise. Both training and test errors remain large because a line cannot capture the curved relationship: the model is too simple for the pattern present.

11.2 Ethics

Use diverse, representative data when possible; report uncertainty; avoid over-claiming; consider impacts on people.

Give one reason a model trained on last year's data may underperform this year.

Hint: Data distribution can shift: the relationship between inputs and outputs changes.

11.3 Further Directions

This book covers feed-forward neural networks—the foundation for more advanced architectures. Transformers, for instance, extend these ideas to handle sequences (text, time series) using attention mechanisms that allow the model to focus on relevant parts of the input when making predictions. Other directions include convolutional networks for images, recurrent networks for sequences, and generative models that learn to create new data.

Remark 11.1

The core concepts you've learned—layers, activations, loss functions, optimisation—apply across these advanced architectures. Understanding the fundamentals prepares you to explore these exciting extensions.

Chapter 12

Glossary

Plain-language definitions for quick recall.

Activation A nonlinear function applied to a neuron's input (e.g., ReLU, sigmoid).

Cost Function Aggregates loss over the entire dataset, usually by averaging; minimised during training. See also *Loss*, *Objective Function*.

Error The difference between a prediction and the true value.

Feature Learning Functions Functions that extract useful representations from data, such as Principal Component Analysis (PCA) for dimensionality reduction and autoencoders for unsupervised feature extraction.

Gradient Descent An iterative method that adjusts parameters to reduce loss.

Kernel Function A function that implicitly maps input data to higher-dimensional spaces for better separability, used in algorithms like Support Vector Machines. The kernel trick avoids explicit computation of transformed features.

Loss A number that measures how far predictions are from targets. See also *Cost Function*, *Objective Function*.

Matrix A rectangular array of numbers; represents a linear transformation.

Neural Network A function built by composing linear steps with activations.

Objective Function A general term encompassing loss and cost functions; represents what the learning algorithm aims to optimise (minimise or maximise). See also *Loss*, *Cost Function*.

Overfitting Learning noise or specifics of training data that do not generalise.

Regularisation A technique that penalises model complexity to prevent overfitting, such as L1 (Lasso) and L2 (Ridge) regularisation added to the cost function.

Vector An ordered list of numbers representing features.

Appendix A

Worksheets: Checkpoints

Short comprehension checks per chapter. Replace this text with specific checkpoint prompts when ready.

Appendix B

Worksheets: Exercises

Small numeric and drawing tasks designed for practice without calculus.

References

Bibliography

- [1] J.S. Cramer. The origins of logistic regression. *Tinbergen Institute Discussion Papers*, (2002-119/4):1–16, 2002.
- [2] Geoffrey E Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [3] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- [4] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [5] Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133, 1943.
- [6] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010.
- [7] Herbert Robbins and Sutton Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- [8] Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386–408, 1958.
- [9] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.
- [10] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

About the Author

Vu Hung Nguyen is the author of this book. For more information and updates:

- Website: <https://vuhung16au.github.io/>
- GitHub: <https://github.com/vuhung16au/>
- LinkedIn: <https://www.linkedin.com/in/nguyenvuhung/>